



















































































































































































































































































































































































































































































































































































































































































































exploratory data analysis and as a preprocessing stage before a supervised learner. One interesting and successful application is the traveling salesman problem (Angeniol, Vaubois, and Le Texier 1988). Just like the difference between  $k$ -means clustering and EM on Gaussian mixtures (chapter 7), *generative topographic mapping* (GTM) (Bishop, Svensén, and Williams 1998) is a probabilistic version of SOM that optimizes the log likelihood of the data using a mixture of Gaussians whose means are constrained to lie on a two-dimensional manifold (for topological ordering in low dimensions).

In an RBF network, once the centers and spreads are fixed (for example, by choosing a random subset of training instances as centers, as in the anchor method), training the second layer is a linear model. This model is equivalent to support vector machines with Gaussian kernels where during learning the best subset of instances, named the *support vectors*, are chosen; we discuss them in chapter 13. Gaussian processes (chapter 14) where we interpolate from stored training instances are also similar.

## 12.11 Exercises

1. Show an RBF network that implements XOR.
2. Write down the RBF network that uses elliptic units instead of radial units as in equation 12.13.
3. Derive the update equations for the RBF network for classification (equations 12.20 and 12.21).
4. Show how the system given in equation 12.22 can be trained.
5. Compare the number of parameters of a mixture of experts architecture with an RBF network.
6. Formalize a mixture of experts architecture where the experts and the gating network are multilayer perceptrons. Derive the update equations for regression and classification.
7. Derive the update equations for the cooperative mixture of experts for classification.
8. Derive the update equations for the competitive mixture of experts for classification.
9. Formalize the hierarchical mixture of experts architecture with two levels. Derive the update equations using gradient descent for regression and classification.

10. In mixture of experts, because different experts specialize in different parts of the input space, they may need to focus on different inputs. Discuss how dimensionality can be locally reduced in the experts.

## 12.12 References

- Alpaydm, E., and M. I. Jordan. 1996. "Local Linear Perceptrons for Classification." *IEEE Transactions on Neural Networks* 7: 788-792.
- Angeniol, B., G. Vaubois, and Y. Le Texier. 1988. "Self Organizing Feature Maps and the Travelling Salesman Problem." *Neural Networks* 1: 289-293.
- Berthold, M. 1999. "Fuzzy Logic." In *Intelligent Data Analysis: An Introduction*, ed. M. Berthold and D. J. Hand, 269-298. Berlin: Springer.
- Bishop, C. M., M. Svensén, and C. K. I. Williams. 1998. "GTM: The Generative Topographic Mapping." *Neural Computation* 10: 215-234.
- Bottou, L., and V. Vapnik. 1992. "Local Learning Algorithms." *Neural Computation* 4: 888-900.
- Broomhead, D. S., and D. Lowe. 1988. "Multivariable Functional Interpolation and Adaptive Networks." *Complex Systems* 2: 321-355.
- Carpenter, G. A., and S. Grossberg. 1988. "The ART of Adaptive Pattern Recognition by a Self-Organizing Neural Network." *IEEE Computer* 21(3): 77-88.
- Cherkassky, V., and F. Mulier. 1998. *Learning from Data: Concepts, Theory, and Methods*. New York: Wiley.
- DeSieno, D. 1988. "Adding a Conscience Mechanism to Competitive Learning." In *IEEE International Conference on Neural Networks*, 117-124. Piscataway, NJ: IEEE Press.
- Feldman, J. A., and D. H. Ballard. 1982. "Connectionist Models and their Properties." *Cognitive Science* 6: 205-254.
- Fritzke, B. 1995. "Growing Cell Structures: A Self Organizing Network for Un-supervised and Supervised Training." *Neural Networks* 7: 1441-1460.
- Grossberg, S. 1980. "How does the Brain Build a Cognitive Code?" *Psychological Review* 87: 1-51.
- Hertz, J., A. Krogh, and R. G. Palmer. 1991. *Introduction to the Theory of Neural Computation*. Reading, MA: Addison Wesley.
- Jacobs, R. A., M. I. Jordan, S. J. Nowlan, and G. E. Hinton. 1991. "Adaptive Mixtures of Local Experts." *Neural Computation* 3: 79-87.
- Jordan, M. I., and R. A. Jacobs. 1994. "Hierarchical Mixtures of Experts and the EM Algorithm." *Neural Computation* 6: 181-214.

- Kohonen, T. 1990. "The Self-Organizing Map." *Proceedings of the IEEE* 78: 1464-1480.
- Kohonen, T. 1995. *Self-Organizing Maps*. Berlin: Springer.
- Lee, Y. 1991. "Handwritten Digit Recognition Using  $k$ -Nearest Neighbor, Radial Basis Function, and Backpropagation Neural Networks." *Neural Computation* 3: 440-449.
- Mao, J., and A. K. Jain. 1995. "Artificial Neural Networks for Feature Extraction and Multivariate Data Projection." *IEEE Transactions on Neural Networks* 6: 296-317.
- Moody, J., and C. Darken. 1989. "Fast Learning in Networks of Locally-Tuned Processing Units." *Neural Computation* 1: 281-294.
- Oja, E. 1982. "A Simplified Neuron Model as a Principal Component Analyzer." *Journal of Mathematical Biology* 15: 267-273.
- Omohundro, S. M. 1987. "Efficient Algorithms with Neural Network Behavior." *Complex Systems* 1: 273-347.
- Platt, J. 1991. "A Resource Allocating Network for Function Interpolation." *Neural Computation* 3: 213-225.
- Specht, D. F. 1991. "A General Regression Neural Network." *IEEE Transactions on Neural Networks* 2: 568-576.
- Tresp, V., J. Hollatz, and S. Ahmad. 1997. "Representing Probabilistic Rules with Networks of Gaussian Basis Functions." *Machine Learning* 27: 173-200.

# 13 *Kernel Machines*

*Kernel machines are maximum margin methods that allow the model to be written as a sum of the influences of a subset of the training instances. These influences are given by application-specific similarity kernels, and we discuss “kernelized” classification, regression, outlier detection, and dimensionality reduction, as well as how to choose and use kernels.*

## 13.1 Introduction

We now discuss a different approach for linear classification and regression. We should not be surprised to have so many different methods even for the simple case of a linear model. Each learning algorithm has a different inductive bias, makes different assumptions, and defines a different objective function and thus may find a different linear model.

The model that we will discuss in this chapter, called the *support vector machine* (SVM), and later generalized under the name *kernel machine*, has been popular in recent years for a number of reasons:

1. It is a discriminant-based method and uses Vapnik’s principle to never solve a more complex problem as a first step before the actual problem (Vapnik 1995). For example, in classification, when the task is to learn the discriminant, it is not necessary to estimate where the class densities  $p(\mathbf{x}|C_i)$  or the exact posterior probability values  $P(C_i|\mathbf{x})$ ; we only need to estimate where the class boundaries lie, that is,  $\mathbf{x}$  where  $P(C_i|\mathbf{x}) = P(C_j|\mathbf{x})$ . Similarly, for outlier detection, we do not need to estimate the full density  $p(\mathbf{x})$ ; we only need to find the boundary separating those  $\mathbf{x}$  that have low  $p(\mathbf{x})$ , that is,  $\mathbf{x}$  where  $p(\mathbf{x}) < \theta$ , for some threshold  $\theta \in (0, 1)$ .

2. After training, the parameter of the linear model, the weight vector, can be written down in terms of a subset of the training set, which are the so-called *support vectors*. In classification, these are the cases that are close to the boundary and as such, knowing them allows knowledge extraction: those are the uncertain or erroneous cases that lie in the vicinity of the boundary between two classes. Their number gives us an estimate of the generalization error, and, as we see below, being able to write the model parameter in terms of a set of instances allows kernelization.
3. As we will see shortly, the output is written as a sum of the influences of support vectors and these are given by *kernel functions* that are application-specific measures of similarity between data instances. Previously, we talked about nonlinear basis functions allowing us to map the input to another space where a linear (smooth) solution is possible; the kernel function uses the same idea.
4. Typically in most learning algorithms, data points are represented as vectors, and either dot product (as in the multilayer perceptrons) or Euclidean distance (as in radial basis function networks) is used. A kernel function allows us to go beyond that. For example,  $G_1$  and  $G_2$  may be two graphs and  $K(G_1, G_2)$  may correspond to the number of shared paths, which we can calculate without needing to represent  $G_1$  or  $G_2$  explicitly as vectors.
5. Kernel-based algorithms are formulated as convex optimization problems, and there is a single optimum that we can solve for analytically. Therefore we are no longer bothered with heuristics for learning rates, initializations, checking for convergence, and such. Of course, this does not mean that we do not have any hyperparameters for model selection; we do—any method needs them, to match the algorithm to the data at hand.

We start our discussion with the case of classification, and then generalize to regression, outlier (novelty) detection, and then dimensionality reduction. We see that in all cases basically we have the similar quadratic program template to maximize the separability, or *margin*, of instances subject to a constraint of the smoothness of solution. Solving for it, we get the support vectors. The kernel function defines the space according to its notion of similarity and a kernel function is good if we have better separation in its corresponding space.

## 13.2 Optimal Separating Hyperplane

Let us start again with two classes and use labels  $-1/ +1$  for the two classes. The sample is  $\mathcal{X} = \{\mathbf{x}^t, r^t\}$  where  $r^t = +1$  if  $\mathbf{x}^t \in C_1$  and  $r^t = -1$  if  $\mathbf{x}^t \in C_2$ . We would like to find  $\mathbf{w}$  and  $w_0$  such that

$$\begin{aligned} \mathbf{w}^T \mathbf{x}^t + w_0 &\geq +1 & \text{for } r^t = +1 \\ \mathbf{w}^T \mathbf{x}^t + w_0 &\leq -1 & \text{for } r^t = -1 \end{aligned}$$

which can be rewritten as

$$(13.1) \quad r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1$$

Note that we do not simply require

$$r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq 0$$

Not only do we want the instances to be on the right side of the hyperplane, but we also want them some distance away, for better generalization. The distance from the hyperplane to the instances closest to it on either side is called the *margin*, which we want to maximize for best generalization.

MARGIN

Very early on, in section 2.1, we talked about the concept of the margin when we were talking about fitting a rectangle, and we said that it is better to take a rectangle halfway between  $S$  and  $G$ , to get a breathing space. This is so that in case noise shifts a test instance slightly, it will still be on the right side of the boundary.

OPTIMAL SEPARATING  
HYPERPLANE

Similarly, now that we are using the hypothesis class of lines, the *optimal separating hyperplane* is the one that maximizes the margin.

We remember from section 10.3 that the distance of  $\mathbf{x}^t$  to the discriminant is

$$\frac{|\mathbf{w}^T \mathbf{x}^t + w_0|}{\|\mathbf{w}\|}$$

which, when  $r^t \in \{-1, +1\}$ , can be written as

$$\frac{r^t (\mathbf{w}^T \mathbf{x}^t + w_0)}{\|\mathbf{w}\|}$$

and we would like this to be at least some value  $\rho$ :

$$(13.2) \quad \frac{r^t (\mathbf{w}^T \mathbf{x}^t + w_0)}{\|\mathbf{w}\|} \geq \rho, \forall t$$

We would like to maximize  $\rho$  but there are an infinite number of solutions that we can get by scaling  $\mathbf{w}$  and for a unique solution, we fix  $\rho\|\mathbf{w}\| = 1$  and thus, to maximize the margin, we minimize  $\|\mathbf{w}\|$ . The task can therefore be defined (see Cortes and Vapnik 1995; Vapnik 1995) as to

$$(13.3) \quad \min \frac{1}{2}\|\mathbf{w}\|^2 \text{ subject to } r^t(\mathbf{w}^T \mathbf{x}^t + w_0) \geq +1, \forall t$$

This is a standard quadratic optimization problem, whose complexity depends on  $d$ , and it can be solved directly to find  $\mathbf{w}$  and  $w_0$ . Then, on both sides of the hyperplane, there will be instances that are  $1/\|\mathbf{w}\|$  away from the hyperplane and the total margin will be  $2/\|\mathbf{w}\|$ .

We saw in section 10.2 that if the problem is not linearly separable, instead of fitting a nonlinear function, one trick is to map the problem to a new space by using nonlinear basis functions. It is generally the case that this new space has many more dimensions than the original space, and, in such a case, we are interested in a method whose complexity does not depend on the input dimensionality.

In finding the optimal hyperplane, we can convert the optimization problem to a form whose complexity depends on  $N$ , the number of training instances, and not on  $d$ . Another advantage of this new formulation is that it will allow us to rewrite the basis functions in terms of kernel functions, as we will see in section 13.5.

To get the new formulation, we first write equation 13.3 as an unconstrained problem using Lagrange multipliers  $\alpha^t$ :

$$(13.4) \quad \begin{aligned} L_p &= \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{t=1}^N \alpha^t [r^t(\mathbf{w}^T \mathbf{x}^t + w_0) - 1] \\ &= \frac{1}{2}\|\mathbf{w}\|^2 - \sum_t \alpha^t r^t(\mathbf{w}^T \mathbf{x}^t + w_0) + \sum_t \alpha^t \end{aligned}$$

This should be minimized with respect to  $\mathbf{w}$ ,  $w_0$  and maximized with respect to  $\alpha^t \geq 0$ . The saddle point gives the solution.

This is a convex quadratic optimization problem because the main term is convex and the linear constraints are also convex. Therefore, we can equivalently solve the dual problem, making use of the Karush-Kuhn-Tucker conditions. The dual is to *maximize*  $L_p$  with respect to  $\alpha^t$ , subject to the constraints that the gradient of  $L_p$  with respect to  $\mathbf{w}$  and  $w_0$  are 0

and also that  $\alpha^t \geq 0$ :

$$(13.5) \quad \frac{\partial L_p}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_t \alpha^t r^t \mathbf{x}^t$$

$$(13.6) \quad \frac{\partial L_p}{\partial w_0} = 0 \Rightarrow \sum_t \alpha^t r^t = 0$$

Plugging these into equation 13.4, we get the dual

$$\begin{aligned} L_d &= \frac{1}{2}(\mathbf{w}^T \mathbf{w}) - \mathbf{w}^T \sum_t \alpha^t r^t \mathbf{x}^t - w_0 \sum_t \alpha^t r^t + \sum_t \alpha^t \\ &= -\frac{1}{2}(\mathbf{w}^T \mathbf{w}) + \sum_t \alpha^t \\ (13.7) \quad &= -\frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\mathbf{x}^t)^T \mathbf{x}^s + \sum_t \alpha^t \end{aligned}$$

which we maximize with respect to  $\alpha^t$  only, subject to the constraints

$$\sum_t \alpha^t r^t = 0, \text{ and } \alpha^t \geq 0, \forall t$$

This can be solved using quadratic optimization methods. The size of the dual depends on  $N$ , sample size, and not on  $d$ , the input dimensionality. The upper bound for time complexity is  $\mathcal{O}(N^3)$ , and the upper bound for space complexity is  $\mathcal{O}(N^2)$ .

Once we solve for  $\alpha^t$ , we see that though there are  $N$  of them, most vanish with  $\alpha^t = 0$  and only a small percentage have  $\alpha^t > 0$ . The set of  $\mathbf{x}^t$  whose  $\alpha^t > 0$  are the *support vectors*, and as we see in equation 13.5,  $\mathbf{w}$  is written as the weighted sum of these training instances that are selected as the support vectors. These are the  $\mathbf{x}^t$  that satisfy

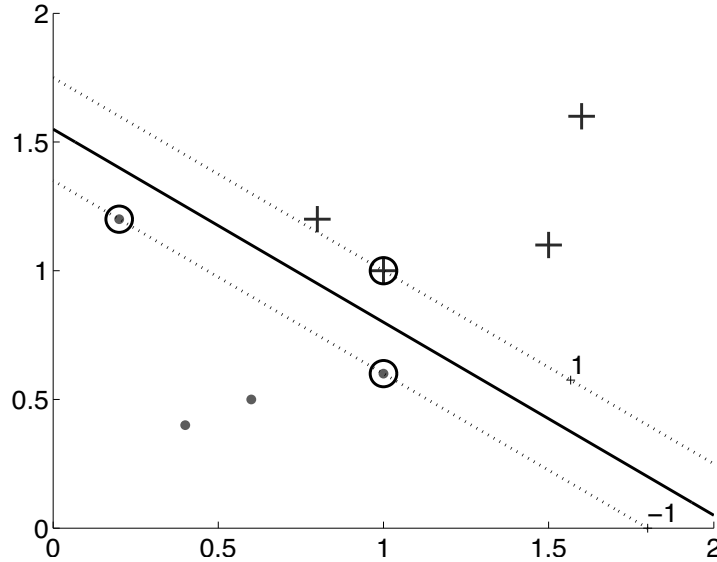
$$r^t (\mathbf{w}^T \mathbf{x}^t + w_0) = 1$$

and lie on the margin. We can use this fact to calculate  $w_0$  from any support vector as

$$(13.8) \quad w_0 = r^t - \mathbf{w}^T \mathbf{x}^t$$

For numerical stability, it is advised that this be done for all support vectors and an average be taken. The discriminant thus found is called the *support vector machine* (SVM) (see figure 13.1).

The majority of the  $\alpha^t$  are 0, for which  $r^t (\mathbf{w}^T \mathbf{x}^t + w_0) > 1$ . These are the  $\mathbf{x}^t$  that lie more than sufficiently away from the discriminant,



**Figure 13.1** For a two-class problem where the instances of the classes are shown by plus signs and dots, the thick line is the boundary and the dashed lines define the margins on either side. Circled instances are the support vectors.

and they have no effect on the hyperplane. The instances that are not support vectors carry no information; even if any subset of them are removed, we would still get the same solution. From this perspective, the SVM algorithm can be likened to the condensed nearest neighbor algorithm (section 8.5), which stores only the instances neighboring (and hence constraining) the class discriminant.

Being a discriminant-based method, the SVM cares only about the instances close to the boundary and discards those that lie in the interior. Using this idea, it is possible to use a simpler classifier before the SVM to filter out a large portion of such instances, thereby decreasing the complexity of the optimization step of the SVM (exercise 1).

During testing, we do not enforce a margin. We calculate  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ , and choose according to the sign of  $g(\mathbf{x})$ :

Choose  $C_1$  if  $g(\mathbf{x}) > 0$  and  $C_2$  otherwise

### 13.3 The Nonseparable Case: Soft Margin Hyperplane

SLACK VARIABLES

If the data is not linearly separable, the algorithm we discussed earlier will not work. In such a case, if the two classes are not linearly separable such that there is no hyperplane to separate them, we look for the one that incurs the least error. We define *slack variables*,  $\xi^t \geq 0$ , which store the deviation from the margin. There are two types of deviation: An instance may lie on the wrong side of the hyperplane and be misclassified. Or, it may be on the right side but may lie in the margin, namely, not sufficiently away from the hyperplane. Relaxing equation 13.1, we require

$$(13.9) \quad r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq 1 - \xi^t$$

If  $\xi^t = 0$ , there is no problem with  $\mathbf{x}^t$ . If  $0 < \xi^t < 1$ ,  $\mathbf{x}^t$  is correctly classified but in the margin. If  $\xi^t \geq 1$ ,  $\mathbf{x}^t$  is misclassified (see figure 13.2). The number of misclassifications is  $\#\{\xi^t > 1\}$ , and the number of non-separable points is  $\#\{\xi^t > 0\}$ . We define *soft error* as

SOFT ERROR

$$\sum_t \xi^t$$

and add this as a penalty term:

$$(13.10) \quad L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi^t$$

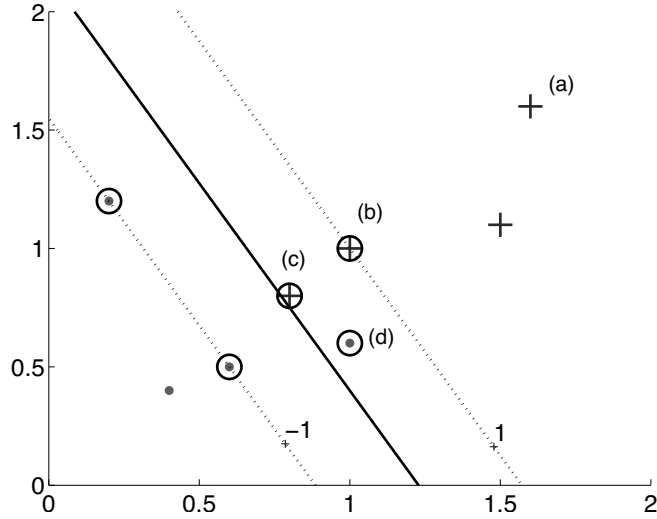
subject to the constraint of equation 13.9.  $C$  is the penalty factor as in any regularization scheme trading off complexity, as measured by the  $L_2$  norm of the weight vector (similar to weight decay in multilayer perceptrons; see section 11.9), and data misfit, as measured by the number of nonseparable points. Note that we are penalizing not only the misclassified points but also the ones in the margin for better generalization, though these latter would be correctly classified during testing.

Adding the constraints, the Lagrangian of equation 13.4 then becomes

$$(13.11) \quad L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi^t - \sum_t \alpha^t [r^t (\mathbf{w}^T \mathbf{x}^t + w_0) - 1 + \xi^t] - \sum_t \mu^t \xi^t$$

where  $\mu_t$  are the new Lagrange parameters to guarantee the positivity of  $\xi^t$ . When we take the derivatives with respect to the parameters and set them to 0, we get:

$$(13.12) \quad \frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} - \sum_t \alpha^t r^t \mathbf{x}^t = 0 \Rightarrow \mathbf{w} = \sum_t \alpha^t r^t \mathbf{x}^t$$



**Figure 13.2** In classifying an instance, there are four possible cases: In (a), the instance is on the correct side and far away from the margin;  $r^t g(\mathbf{x}^t) > 1$ ,  $\xi^t = 0$ . In (b),  $\xi^t = 0$ ; it is on the right side and on the margin. In (c),  $\xi^t = 1 - g(\mathbf{x}^t)$ ,  $0 < \xi < 1$ ; it is on the right side but is in the margin and not sufficiently away. In (d),  $\xi^t = 1 + g(\mathbf{x}^t) > 1$ ; it is on the wrong side—this is a misclassification. All cases except (a) are support vectors. In terms of the dual variable, in (a),  $\alpha^t = 0$ ; in (b),  $\alpha^t < C$ ; in (c) and (d),  $\alpha^t = C$ .

$$(13.13) \quad \frac{\partial L_p}{\partial w_0} = \sum_t \alpha^t r^t = 0$$

$$(13.14) \quad \frac{\partial L_p}{\partial \xi^t} = C - \alpha^t - \mu^t = 0$$

Since  $\mu^t \geq 0$ , this last implies that  $0 \leq \alpha^t \leq C$ . Plugging these into equation 13.11, we get the dual that we maximize with respect to  $\alpha^t$ :

$$(13.15) \quad L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\mathbf{x}^t)^T \mathbf{x}^s$$

subject to

$$\sum_t \alpha^t r^t = 0 \text{ and } 0 \leq \alpha^t \leq C, \forall t$$

Solving this, we see that as in the separable case, instances that lie on the correct side of the boundary with sufficient margin vanish with their  $\alpha^t = 0$  (see figure 13.2). The support vectors have their  $\alpha^t > 0$  and they define  $\mathbf{w}$ , as given in equation 13.12. Of these, those whose  $\alpha^t < C$  are the ones that are on the margin, and we can use them to calculate  $w_0$ ; they have  $\xi^t = 0$  and satisfy  $r^t(\mathbf{w}^T \mathbf{x}^t + w_0) = 1$ . Again, it is better to take an average over these  $w_0$  estimates. Those instances that are in the margin or misclassified have their  $\alpha^t = C$ .

The nonseparable instances that we store as support vectors are the instances that we would have trouble correctly classifying if they were not in the training set; they would either be misclassified or classified correctly but not with enough confidence. We can say that the number of support vectors is an upper-bound estimate for the expected number of errors. And, actually, Vapnik (1995) has shown that the expected test error rate is

$$E_N[P(\text{error})] \leq \frac{E_N[\# \text{ of support vectors}]}{N}$$

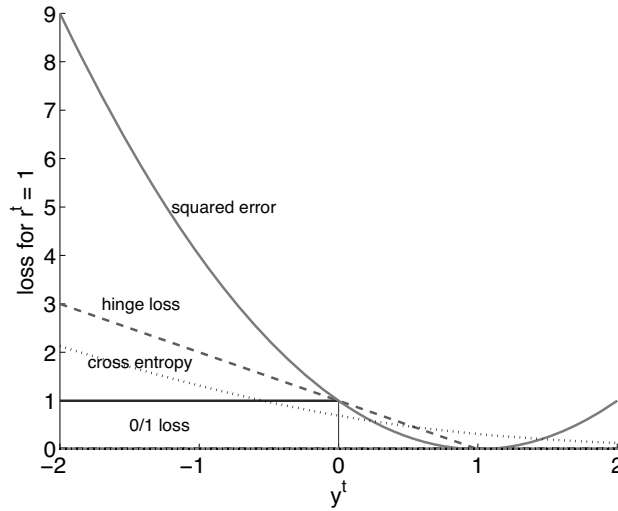
where  $E_N[\cdot]$  denotes expectation over training sets of size  $N$ . The nice implication of this is that it shows that the error rate depends on the number of support vectors and not on the input dimensionality.

Equation 13.9 implies that we define error if the instance is on the wrong side or if the margin is less than 1. This is called the *hinge loss*. If  $y^t = \mathbf{w}^T \mathbf{x}^t + w_0$  is the output and  $r^t$  is the desired output, hinge loss is defined as

$$(13.16) \quad L_{\text{hinge}}(y^t, r^t) = \begin{cases} 0 & \text{if } y^t r^t \geq 1 \\ 1 - y^t r^t & \text{otherwise} \end{cases}$$

In figure 13.3, we compare hinge loss with 0/1 loss, squared error, and cross-entropy. We see that different from 0/1 loss, hinge loss also penalizes instances in the margin even though they may be on the correct side, and the loss increases linearly as the instance moves away on the wrong side. This is different from the squared loss that therefore is not as robust as the hinge loss. We see that cross-entropy minimized in logistic discrimination (section 10.7) or by the linear perceptron (section 11.3), is a good continuous approximation to the hinge loss.

$C$  of equation 13.10 is the regularization parameter fine-tuned using cross-validation. It defines the trade-off between margin maximization and error minimization: If it is too large, we have a high penalty for nonseparable points, and we may store many support vectors and overfit.



**Figure 13.3** Comparison of different loss functions for  $r^t = 1$ : 0/1 loss is 0 if  $y^t = 1$ , 1 otherwise. Hinge loss is 0 if  $y^t > 1$ ,  $1 - y^t$  otherwise. Squared error is  $(1 - y^t)^2$ . Cross-entropy is  $\log(1/(1 + \exp(-y^t)))$ .

If it is too small, we may find too simple solutions that underfit. Typically, one chooses from  $[10^{-6}, 10^{-5}, \dots, 10^{+5}, 10^{+6}]$  in the log scale by looking at the accuracy on a validation set.

### 13.4 $\nu$ -SVM

There is another, equivalent formulation of the soft margin hyperplane that uses a parameter  $\nu \in [0, 1]$  instead of  $C$  (Schölkopf et al. 2000). The objective function is

$$(13.17) \quad \min \frac{1}{2} \|\mathbf{w}\|^2 - \nu \rho + \frac{1}{N} \sum_t \xi^t$$

subject to

$$(13.18) \quad r^t (\mathbf{w}^T \mathbf{x}^t + w_0) \geq \rho - \xi^t, \quad \xi^t \geq 0, \quad \rho \geq 0$$

$\rho$  is a new parameter that is a variable of the optimization problem and scales the margin: the margin is now  $2\rho/\|\mathbf{w}\|$ .  $\nu$  has been shown to be a

lower bound on the fraction of support vectors and an upper bound on the fraction of instances having margin errors ( $\sum_t \#\{\xi^t > 0\}$ ). The dual is

$$(13.19) \quad L_d = -\frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s (\mathbf{x}^t)^T \mathbf{x}^s$$

subject to

$$\sum_t \alpha^t r^t = 0, \quad 0 \leq \alpha^t \leq \frac{1}{N}, \quad \sum_t \alpha^t \leq \nu$$

When we compare equation 13.19 with equation 13.15, we see that the term  $\sum_t \alpha^t$  no longer appears in the objective function but is now a constraint. By playing with  $\nu$ , we can control the fraction of support vectors, and this is advocated to be more intuitive than playing with  $C$ .

### 13.5 Kernel Trick

Section 10.2 demonstrated that if the problem is nonlinear, instead of trying to fit a nonlinear model, we can map the problem to a new space by doing a nonlinear transformation using suitably chosen basis functions and then use a linear model in this new space. The linear model in the new space corresponds to a nonlinear model in the original space. This approach can be used in both classification and regression problems, and in the special case of classification, it can be used with any scheme. In the particular case of support vector machines, it leads to certain simplifications that we now discuss.

Let us say we have the new dimensions calculated through the basis functions

$$\mathbf{z} = \boldsymbol{\phi}(\mathbf{x}) \text{ where } z_j = \phi_j(\mathbf{x}), j = 1, \dots, k$$

mapping from the  $d$ -dimensional  $\mathbf{x}$  space to the  $k$ -dimensional  $\mathbf{z}$  space where we write the discriminant as

$$(13.20) \quad \begin{aligned} g(\mathbf{z}) &= \mathbf{w}^T \mathbf{z} \\ g(\mathbf{x}) &= \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \\ &= \sum_{j=1}^k w_j \phi_j(\mathbf{x}) \end{aligned}$$

where we do not use a separate  $w_0$ ; we assume that  $z_1 = \phi_1(\mathbf{x}) \equiv 1$ . Generally,  $k$  is much larger than  $d$  and  $k$  may also be larger than  $N$ , and there

lies the advantage of using the dual form whose complexity depends on  $N$ , whereas if we used the primal it would depend on  $k$ . We also use the more general case of the soft margin hyperplane here because we have no guarantee that the problem is linearly separable in this new space.

The problem is the same

$$(13.21) \quad L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi^t$$

except that now the constraints are defined in the new space

$$(13.22) \quad r^t \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^t) \geq 1 - \xi^t$$

The Lagrangian is

$$(13.23) \quad L_p = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t \xi^t - \sum_t \alpha^t [r^t \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^t) - 1 + \xi^t] - \sum_t \mu^t \xi^t$$

When we take the derivatives with respect to the parameters and set them to 0, we get

$$(13.24) \quad \frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} = \sum_t \alpha^t r^t \boldsymbol{\phi}(\mathbf{x}^t)$$

$$(13.25) \quad \frac{\partial L_p}{\partial \xi^t} = C - \alpha^t - \mu^t = 0$$

The dual is now

$$(13.26) \quad L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s \boldsymbol{\phi}(\mathbf{x}^t)^T \boldsymbol{\phi}(\mathbf{x}^s)$$

subject to

$$\sum_t \alpha^t r^t = 0 \text{ and } 0 \leq \alpha^t \leq C, \forall t$$

KERNEL FUNCTION

The idea in *kernel machines* is to replace the inner product of basis functions,  $\boldsymbol{\phi}(\mathbf{x}^t)^T \boldsymbol{\phi}(\mathbf{x}^s)$ , by a *kernel function*,  $K(\mathbf{x}^t, \mathbf{x}^s)$ , between instances in the original input space. So instead of mapping two instances  $\mathbf{x}^t$  and  $\mathbf{x}^s$  to the  $\mathbf{z}$ -space and doing a dot product there, we directly apply the kernel function in the original space.

$$(13.27) \quad L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s K(\mathbf{x}^t, \mathbf{x}^s)$$

The kernel function also shows up in the discriminant

$$\begin{aligned}
 g(\mathbf{x}) &= \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) = \sum_t \alpha^t r^t \boldsymbol{\phi}(\mathbf{x}^t)^T \boldsymbol{\phi}(\mathbf{x}) \\
 (13.28) \quad &= \sum_t \alpha^t r^t K(\mathbf{x}^t, \mathbf{x})
 \end{aligned}$$

This implies that if we have the kernel function, we do not need to map it to the new space at all. Actually, for any valid kernel, there does exist a corresponding mapping function, but it may be much simpler to use  $K(\mathbf{x}^t, \mathbf{x})$  rather than calculating  $\boldsymbol{\phi}(\mathbf{x}^t)$ ,  $\boldsymbol{\phi}(\mathbf{x})$  and taking the dot product. Many algorithms have been *kernelized*, as we will see in later sections, and that is why we have the name “kernel machines.”

KERNELIZATION

GRAM MATRIX

The matrix of kernel values,  $\mathbf{K}$ , where  $\mathbf{K}_{ts} = K(\mathbf{x}^t, \mathbf{x}^s)$ , is called the *Gram matrix*, which should be symmetric and positive semidefinite. Recently, it has become standard practice in sharing data sets to have available only the  $\mathbf{K}$  matrices without providing  $\mathbf{x}^t$  or  $\boldsymbol{\phi}(\mathbf{x}^t)$ . Especially in bioinformatics or natural language processing applications where  $\mathbf{x}$  (or  $\boldsymbol{\phi}(\mathbf{x})$ ) has hundreds or thousands of dimensions, storing/downloading the  $N \times N$  matrix is much cheaper (Vert, Tsuda, and Schölkopf 2004); this, however, implies that we can use only those available for training/testing and cannot use the trained model to make predictions outside this data set.

## 13.6 Vectorial Kernels

The most popular, general-purpose kernel functions are

- *polynomials* of degree  $q$ :

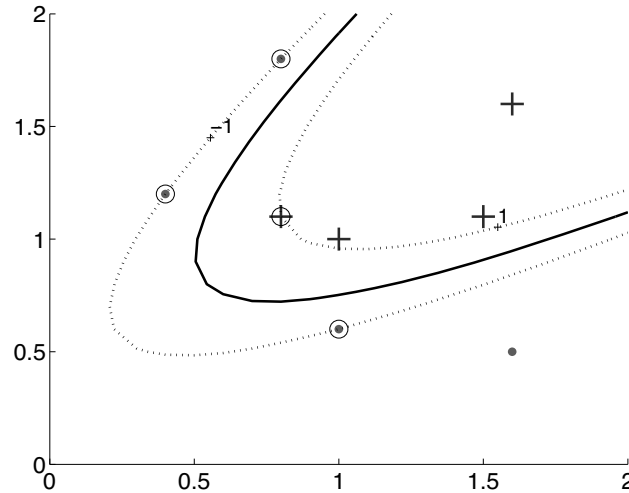
$$(13.29) \quad K(\mathbf{x}^t, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t + 1)^q$$

where  $q$  is selected by the user. For example, when  $q = 2$  and  $d = 2$ ,

$$\begin{aligned}
 K(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^T \mathbf{y} + 1)^2 \\
 &= (x_1 y_1 + x_2 y_2 + 1)^2 \\
 &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2
 \end{aligned}$$

corresponds to the inner product of the basis function (Cherkassky and Mulier 1998):

$$\boldsymbol{\phi}(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2]^T$$



**Figure 13.4** The discriminant and margins found by a polynomial kernel of degree 2. Circled instances are the support vectors.

An example is given in figure 13.4. When  $q = 1$ , we have the *linear kernel* that corresponds to the original formulation.

■ *radial-basis functions:*

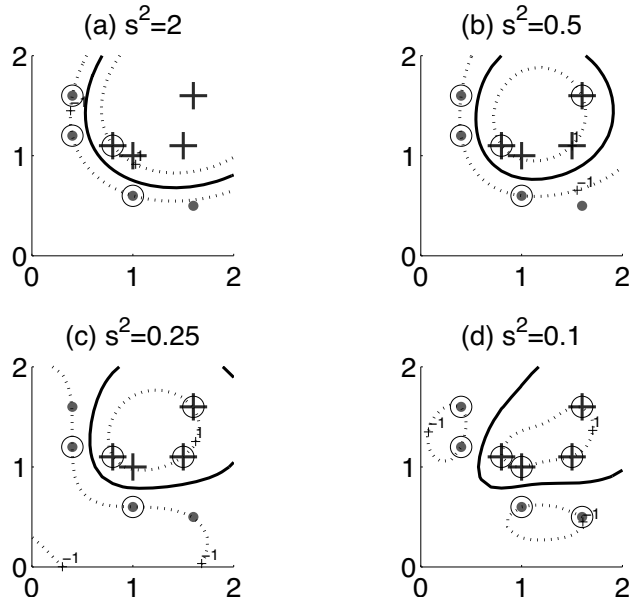
$$(13.30) \quad K(\mathbf{x}^t, \mathbf{x}) = \exp \left[ -\frac{\|\mathbf{x}^t - \mathbf{x}\|^2}{2s^2} \right]$$

defines a spherical kernel as in Parzen windows (chapter 8) where  $\mathbf{x}^t$  is the center and  $s$ , supplied by the user, defines the radius. This is also similar to radial basis functions that we discuss in chapter 12.

An example is shown in figure 13.5 where we see that larger spreads smooth the boundary; the best value is found by cross-validation. Note that when there are two parameters to be optimized using cross-validation, for example, here  $C$  and  $s^2$ , one should do a grid (factorial) search in the two dimensions; we will discuss methods for searching the best combination of such factors in section 19.2.

One can have a Mahalanobis kernel, generalizing from the Euclidean distance:

$$(13.31) \quad K(\mathbf{x}^t, \mathbf{x}) = \exp \left[ -\frac{1}{2}(\mathbf{x}^t - \mathbf{x})^T \mathbf{S}^{-1}(\mathbf{x}^t - \mathbf{x}) \right]$$



**Figure 13.5** The boundary and margins found by the Gaussian kernel with different spread values,  $s^2$ . We get smoother boundaries with larger spreads.

where  $\mathbf{S}$  is a covariance matrix. Or, in the most general case,

$$(13.32) \quad K(\mathbf{x}^t, \mathbf{x}) = \exp \left[ -\frac{\mathcal{D}(\mathbf{x}^t, \mathbf{x})}{2s^2} \right]$$

for some distance function  $\mathcal{D}(\mathbf{x}^t, \mathbf{x})$ .

■ *sigmoidal functions:*

$$(13.33) \quad K(\mathbf{x}^t, \mathbf{x}) = \tanh(2\mathbf{x}^T \mathbf{x}^t + 1)$$

where  $\tanh(\cdot)$  has the same shape with sigmoid, except that it ranges between  $-1$  and  $+1$ . This is similar to multilayer perceptrons that we discussed in chapter 11.

## 13.7 Defining Kernels

It is also possible to define application-specific kernels. Kernels are generally considered to be measures of similarity in the sense that  $K(\mathbf{x}, \mathbf{y})$  takes a larger value as  $\mathbf{x}$  and  $\mathbf{y}$  are more “similar,” from the point of view of the application. This implies that any prior knowledge we have regarding the application can be provided to the learner through appropriately defined kernels—“kernel engineering”—and such use of kernels can be seen as another example of a “hint” (section 11.8.4).

There are string kernels, tree kernels, graph kernels, and so on (Vert, Tsuda, and Schölkopf 2004), depending on how we represent the data and how we measure similarity in that representation.

BAG OF WORDS

For example, given two documents, the number of words appearing in both may be a kernel. Let us say  $D_1$  and  $D_2$  are two documents and one possible representation is called *bag of words* where we predefine  $M$  words relevant for the application, and we define  $\boldsymbol{\phi}(D_1)$  as the  $M$ -dimensional binary vector whose dimension  $i$  is 1 if word  $i$  appears in  $D_1$  and is 0 otherwise. Then,  $\boldsymbol{\phi}(D_1)^T \boldsymbol{\phi}(D_2)$  counts the number of shared words. Here, we see that if we directly define and implement  $K(D_1, D_2)$  as the number of shared words, we do not need to preselect  $M$  words and can use just any word in the vocabulary (of course, after discarding uninformative words like “of,” “and,” etc.) and we would not need to generate the bag-of-words representation explicitly and it would be as if we allowed  $M$  to be as large as we want.

EDIT DISTANCE

Sometimes—for example, in bioinformatics applications—we can calculate a *similarity score* between two objects, which may not necessarily be positive semidefinite. Given two strings (of genes), a kernel measures the *edit distance*, namely, how many operations (insertions, deletions, substitutions) it takes to convert one string into another; this is also called *alignment*. In such a case, a trick is to define a set of  $M$  templates and represent an object as the  $M$ -dimensional vector of scores to all the templates. That is, if  $\mathbf{m}_i, i = 1, \dots, M$  are the templates and  $s(\mathbf{x}^t, \mathbf{m}_i)$  is the score between  $\mathbf{x}^t$  and  $\mathbf{m}_i$ , then we define

ALIGNMENT

$$\boldsymbol{\phi}(\mathbf{x}^t) = [s(\mathbf{x}^t, \mathbf{m}_1), s(\mathbf{x}^t, \mathbf{m}_2), \dots, s(\mathbf{x}^t, \mathbf{m}_M)]^T$$

EMPIRICAL KERNEL  
MAP

and we define the *empirical kernel map* as

$$K(\mathbf{x}^t, \mathbf{x}^s) = \boldsymbol{\phi}(\mathbf{x}^t)^T \boldsymbol{\phi}(\mathbf{x}^s)$$

which is a valid kernel.

Sometimes, we have a binary score function; for example, two proteins may interact or not, and we want to be able to generalize from this to scores for two arbitrary instances. In such a case, a trick is to define a graph where the nodes are the instances and two nodes are linked if they interact, that is, if the binary score returns 1. Then we say that two nodes that are not immediately linked are “similar” if the path between them is short or if they are connected by many paths. This converts pairwise local interactions to a global similarity measure, rather like defining a geodesic distance used in Isomap (section 6.7), and it is called the *diffusion kernel*.

DIFFUSION KERNEL

If  $p(\mathbf{x})$  is a probability density, then

$$K(\mathbf{x}^t, \mathbf{x}) = p(\mathbf{x}^t)p(\mathbf{x})$$

is a valid kernel. This is used when  $p(\mathbf{x})$  is a generative model for  $\mathbf{x}$  measuring how likely it is that we see  $\mathbf{x}$ . For example, if  $\mathbf{x}$  is a sequence,  $p(\mathbf{x})$  can be a hidden Markov model (chapter 15). With this kernel,  $K(\mathbf{x}^t, \mathbf{x})$  will take a high value if both  $\mathbf{x}^t$  and  $\mathbf{x}$  are likely to have been generated by the same model. It is also possible to parametrize the generative model as  $p(\mathbf{x}|\theta)$  and learn  $\theta$  from data; this is called the *Fisher kernel* (Jaakkola and Haussler 1998).

FISHER KERNEL

## 13.8 Multiple Kernel Learning

It is possible to construct new kernels by combining simpler kernels. If  $K_1(\mathbf{x}, \mathbf{y})$  and  $K_2(\mathbf{x}, \mathbf{y})$  are valid kernels and  $c$  a constant, then

$$(13.34) \quad K(\mathbf{x}, \mathbf{y}) = \begin{cases} cK_1(\mathbf{x}, \mathbf{y}) \\ K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y}) \\ K_1(\mathbf{x}, \mathbf{y}) \cdot K_2(\mathbf{x}, \mathbf{y}) \end{cases}$$

are also valid.

Different kernels may also be using different subsets of  $\mathbf{x}$ . We can therefore see combining kernels as another way to fuse information from different sources where each kernel measures similarity according to its domain. When we have input from two representations  $A$  and  $B$

$$(13.35) \quad \begin{aligned} K_A(\mathbf{x}_A, \mathbf{y}_A) + K_B(\mathbf{x}_B, \mathbf{y}_B) &= \boldsymbol{\phi}_A(\mathbf{x}_A)^T \boldsymbol{\phi}_A(\mathbf{y}_A) + \boldsymbol{\phi}_B(\mathbf{x}_B)^T \boldsymbol{\phi}_B(\mathbf{y}_B) \\ &= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{y}) \\ &= K(\mathbf{x}, \mathbf{y}) \end{aligned}$$

where  $\mathbf{x} = [\mathbf{x}_A, \mathbf{x}_B]$  is the concatenation of the two representations. That is, taking a sum of two kernels corresponds to doing a dot product in the concatenated feature vectors. One can generalize to a number of kernels

$$(13.36) \quad K(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m K_i(\mathbf{x}, \mathbf{y})$$

which, similar to taking an average of classifiers (section 17.4), this time averages over kernels and frees us from the need to choose one particular kernel. It is also possible to take a weighted sum and also learn the weights from data (Lanckriet et al. 2004; Sonnenburg et al. 2006):

$$(13.37) \quad K(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^m \eta_i K_i(\mathbf{x}, \mathbf{y})$$

MULTIPLE KERNEL  
LEARNING

subject to  $\eta_i \geq 0$ , with or without the constraint of  $\sum_i \eta_i = 1$ , respectively known as convex or conic combination. This is called *multiple kernel learning* where we replace a single kernel with a weighted sum. The single kernel objective function of equation 13.27 becomes

$$(13.38) \quad L_d = \sum_t \alpha^t - \frac{1}{2} \sum_t \sum_s \alpha^t \alpha^s r^t r^s \sum_i \eta_i K_i(\mathbf{x}^t, \mathbf{x}^s)$$

which we solve for both the support vector machine parameters  $\alpha^t$  and the kernel weights  $\eta_i$ . Then, the combination of multiple kernels also appear in the discriminant

$$(13.39) \quad g(\mathbf{x}) = \sum_t \alpha^t r^t \sum_i \eta_i K_i(\mathbf{x}^t, \mathbf{x})$$

After training,  $\eta_i$  will take values depending on how the corresponding kernel  $K_i(\mathbf{x}^t, \mathbf{x})$  is useful in discriminating. It is also possible to localize kernels by defining kernel weights as a parameterized function of the input  $\mathbf{x}$ , rather like the gating function in mixture of experts (section 17.8)

$$(13.40) \quad g(\mathbf{x}) = \sum_t \alpha^t r^t \sum_i \eta_i(\mathbf{x}|\theta_i) K_i(\mathbf{x}^t, \mathbf{x})$$

and the gating parameters  $\theta_i$  are learned together with the support vector machine parameters (Gönen and Alpaydm 2008).

When we have information coming from multiple sources in different representations or modalities—for example, in speech recognition where we may have both acoustic and visual lip image—the usual approach is to feed them separately to different classifiers and then fuse the decisions;

we will discuss methods for this in detail in chapter 17. Combining multiple kernels provides us with another way of integrating input from multiple sources, where there is a single classifier that uses different kernels for inputs of different sources, for which there are different notions of similarity (Noble 2004). The localized version can then be seen as an extension of this where we can choose between sources, and hence similarity measures, depending on the input.

### 13.9 Multiclass Kernel Machines

When there are  $K > 2$  classes, the straightforward, *one-vs.-all* way is to define  $K$  two-class problems, each one separating one class from all other classes combined and learn  $K$  support vector machines  $g_i(\mathbf{x}), i = 1, \dots, K$ . That is, in training  $g_i(\mathbf{x})$ , examples of  $C_i$  are labeled +1 and examples of  $C_k, k \neq i$  are labeled as -1. During testing, we calculate all  $g_i(\mathbf{x})$  and choose the maximum.

Platt (1999) proposed to fit a sigmoid to the output of a single (2-class) SVM output to convert to a posterior probability. Similarly, one can train one layer of softmax outputs to minimize cross-entropy to generate  $K > 2$  posterior probabilities (Mayraz and Alpaydm 1999):

$$(13.41) \quad y_i(\mathbf{x}) = \sum_{j=1}^K v_{ij} f_j(\mathbf{x}) + v_{i0}$$

where  $f_j(\mathbf{x})$  are the SVM outputs and  $y_i$  are the posterior probability outputs. Weights  $v_{ij}$  are trained to minimize cross-entropy. Note, however, that as in stacking (section 17.9), the data on which we train  $v_{ij}$  should be different from the data used to train the base SVMs  $f_j(\mathbf{x})$ , to alleviate overfitting.

Instead of the usual approach of building  $K$  two-class SVM classifiers to separate one from all the rest, as with any other classifier, one can build  $K(K-1)/2$  *pairwise* classifiers (see also section 10.4), each  $g_{ij}(\mathbf{x})$  taking examples of  $C_i$  with the label +1, examples of  $C_j$  with the label -1, and not using examples of the other classes. Separating classes in pairs is normally expected to be an easier job, with the additional advantage that because we use less data, the optimizations will be faster, noting however that we have  $\mathcal{O}(K^2)$  discriminants to train instead of  $\mathcal{O}(K)$ .

In the general case, both one-vs.-all and pairwise separation are special cases of the *error-correcting output codes* that decompose a multiclass

ECOC

problem to a set of two-class problems (Dietterich and Bakiri 1995) (see also section 17.6). SVMs being two-class classifiers are ideally suited to this (Allwein, Schapire, and Singer 2000), and it is also possible to have an incremental approach where new two-class SVMs are added to better separate pairs of classes that are confused, to ameliorate a poor ECOC matrix (Mayoraz and Alpaydm 1999).

Another possibility is to write a single *multiclass* optimization problem involving all classes (Weston and Watkins 1998):

$$(13.42) \quad \min \frac{1}{2} \sum_{i=1}^K \|\mathbf{w}_i\|^2 + C \sum_i \sum_t \xi_i^t$$

subject to

$$\mathbf{w}_{z^t} \mathbf{x}^t + w_{z^t 0} \geq \mathbf{w}_i \mathbf{x}^t + w_{i0} + 2 - \xi_i^t, \forall i \neq z^t \text{ and } \xi_i^t \geq 0$$

where  $z^t$  contains the class index of  $\mathbf{x}^t$ . The regularization terms minimize the norms of all hyperplanes simultaneously, and the constraints are there to make sure that the margin between the actual class and any other class is at least 2. The output for the correct class should be at least +1, the output of any other class should be at least -1, and the slack variables are defined to make up any difference.

Though this looks neat, the one-vs.-all approach is generally preferred because it solves  $K$  separate  $N$  variable problems whereas the multiclass formulation uses  $K \cdot N$  variables.

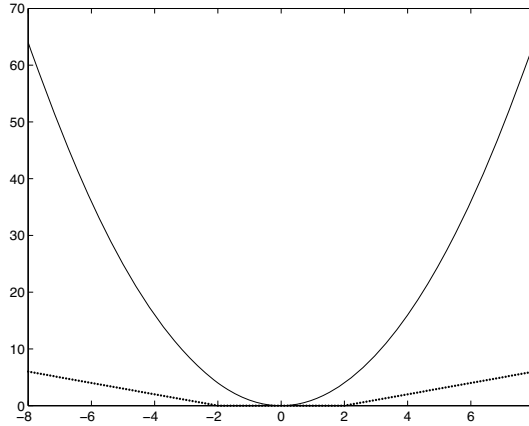
### 13.10 Kernel Machines for Regression

Now let us see how support vector machines can be generalized for regression. We see that the same approach of defining acceptable margins, slacks, and a regularizing function that combines smoothness and error is also applicable here. We start with a linear model, and later on we see how we can use kernel functions here as well:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

In regression proper, we use the square of the difference as error:

$$e_2(r^t, f(\mathbf{x}^t)) = [r^t - f(\mathbf{x}^t)]^2$$



**Figure 13.6** Quadratic and  $\epsilon$ -sensitive error functions. We see that  $\epsilon$ -sensitive error function is not affected by small errors and also is affected less by large errors and thus is more robust to outliers.

whereas in support vector regression, we use the  $\epsilon$ -sensitive loss function:

$$(13.43) \quad e_\epsilon(r^t, f(\mathbf{x}^t)) = \begin{cases} 0 & \text{if } |r^t - f(\mathbf{x}^t)| < \epsilon \\ |r^t - f(\mathbf{x}^t)| - \epsilon & \text{otherwise} \end{cases}$$

ROBUST REGRESSION

which means that we tolerate errors up to  $\epsilon$  and also that errors beyond have a linear effect and not a quadratic one. This error function is therefore more tolerant to noise and is thus more *robust* (see figure 13.6). As in the hinge loss, there is a region of no error, which causes sparseness.

Analogous to the soft margin hyperplane, we introduce slack variables to account for deviations out of the  $\epsilon$ -zone and we get (Vapnik 1995)

$$(13.44) \quad \min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t (\xi_+^t + \xi_-^t)$$

subject to

$$\begin{aligned} r^t - (\mathbf{w}^T \mathbf{x} + w_0) &\leq \epsilon + \xi_+^t \\ (\mathbf{w}^T \mathbf{x} + w_0) - r^t &\leq \epsilon + \xi_-^t \\ \xi_+^t, \xi_-^t &\geq 0 \end{aligned}$$

where we use two types of slack variables, for positive and negative deviations, to keep them positive. Actually, we can see this as two hinges

added back to back, one for positive and one for negative slacks. This formulation corresponds to the  $\epsilon$ -sensitive loss function given in equation 13.43. The Lagrangian is

$$\begin{aligned}
 L_p &= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_t (\xi_+^t + \xi_-^t) \\
 &\quad - \sum_t \alpha_+^t [\epsilon + \xi_+^t - r^t + (\mathbf{w}^T \mathbf{x} + w_0)] \\
 &\quad - \sum_t \alpha_-^t [\epsilon + \xi_-^t + (\mathbf{w}^T \mathbf{x} + w_0) - r^t] \\
 (13.45) \quad &\quad - \sum_t (\mu_+^t \xi_+^t + \mu_-^t \xi_-^t)
 \end{aligned}$$

Taking the partial derivatives, we get

$$(13.46) \quad \frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} - \sum_t (\alpha_+^t - \alpha_-^t) \mathbf{x}^t = 0 \Rightarrow \mathbf{w} = \sum_t (\alpha_+^t - \alpha_-^t) \mathbf{x}^t$$

$$(13.47) \quad \frac{\partial L_p}{\partial w_0} = \sum_t (\alpha_+^t - \alpha_-^t) \mathbf{x}^t = 0$$

$$(13.48) \quad \frac{\partial L_p}{\partial \xi_+^t} = C - \alpha_+^t - \mu_+^t = 0$$

$$(13.49) \quad \frac{\partial L_p}{\partial \xi_-^t} = C - \alpha_-^t - \mu_-^t = 0$$

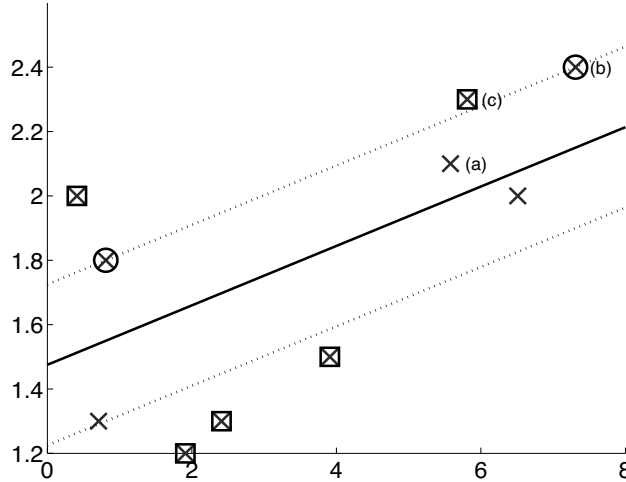
The dual is

$$\begin{aligned}
 L_d &= -\frac{1}{2} \sum_t \sum_s (\alpha_+^t - \alpha_-^t) (\alpha_+^s - \alpha_-^s) (\mathbf{x}^t)^T \mathbf{x}^s \\
 (13.50) \quad &\quad - \epsilon \sum_t (\alpha_+^t + \alpha_-^t) - \sum_t r^t (\alpha_+^t - \alpha_-^t)
 \end{aligned}$$

subject to

$$0 \leq \alpha_+^t \leq C, 0 \leq \alpha_-^t \leq C, \sum_t (\alpha_+^t - \alpha_-^t) = 0$$

Once we solve this, we see that all instances that fall in the tube have  $\alpha_+^t = \alpha_-^t = 0$ ; these are the instances that are fitted with enough precision (see figure 13.7). The support vectors satisfy either  $\alpha_+^t > 0$  or  $\alpha_-^t > 0$  and are of two types. They may be instances that are on the boundary of the tube (either  $\alpha_+^t$  or  $\alpha_-^t$  is between 0 and  $C$ ), and we use these to calculate  $w_0$ . For example, assuming that  $\alpha_+^t > 0$ , we have  $r^t = \mathbf{x}^T \mathbf{x}^t + w_0 + \epsilon$ . Instances that fall outside the  $\epsilon$ -tube are of the second type; these are



**Figure 13.7** The fitted regression line to data points shown as crosses and the  $\epsilon$ -tube are shown ( $C = 10, \epsilon = 0.25$ ). There are three cases: In (a), the instance is in the tube; in (b), the instance is on the boundary of the tube (circled instances); in (c), it is outside the tube with a positive slack, that is,  $\xi_+^t > 0$  (squared instances). (b) and (c) are support vectors. In terms of the dual variable, in (a),  $\alpha_+^t = 0, \alpha_-^t = 0$ , in (b),  $\alpha_+^t < C$ , and in (c),  $\alpha_+^t = C$ .

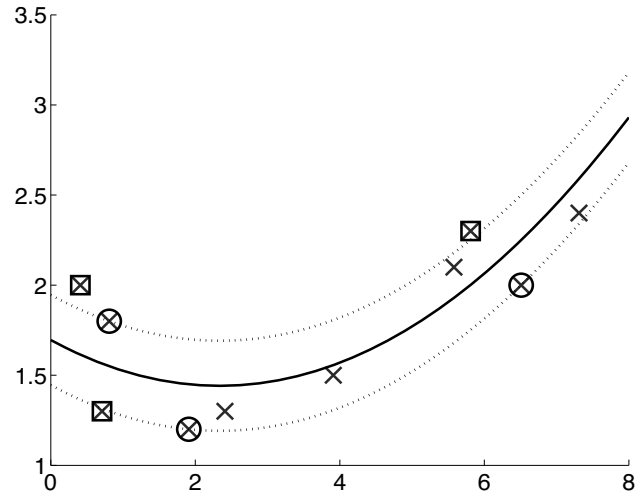
instances for which we do not have a good fit ( $\alpha_+^t = C$ ), as shown in figure 13.7.

Using equation 13.46, we can write the fitted line as a weighted sum of the support vectors:

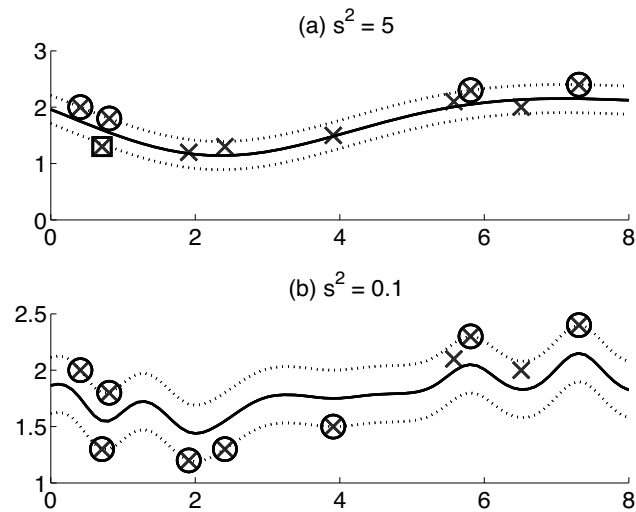
$$(13.51) \quad f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \sum_t (\alpha_+^t - \alpha_-^t) (\mathbf{x}^t)^T \mathbf{x} + w_0$$

Again, the dot product  $(\mathbf{x}^t)^T \mathbf{x}^s$  in equation 13.50 can be replaced with a kernel  $K(\mathbf{x}^t, \mathbf{x}^s)$ , and similarly  $(\mathbf{x}^t)^T \mathbf{x}$  be replaced with  $K(\mathbf{x}^t, \mathbf{x})$  and we can have a nonlinear fit. Using a polynomial kernel would be similar to fitting a polynomial (figure 13.8), and using a Gaussian kernel (figure 13.9) would be similar to nonparametric smoothing models (section 8.6) except that because of the sparsity of solution, we would not need the whole training set but only a subset.

There is also an equivalent  $\nu$ -SVM formulation for regression (Schölkopf et al. 2000), where instead of fixing  $\epsilon$ , we fix  $\nu$  to bound the fraction of support vectors. There is still a need for  $C$  though.



**Figure 13.8** The fitted regression line and the  $\epsilon$ -tube using a quadratic kernel are shown ( $C = 10, \epsilon = 0.25$ ). Circled instances are the support vectors on the margins, squared instances are support vectors which are outliers.



**Figure 13.9** The fitted regression line and the  $\epsilon$ -tube using a Gaussian kernel with two different spreads are shown ( $C = 10, \epsilon = 0.25$ ). Circled instances are the support vectors on the margins, and squared instances are support vectors that are outliers.

### 13.11 One-Class Kernel Machines

Support vector machines, originally proposed for classification, are extended to regression by defining slack variables for deviations around the regression line, instead of the discriminant. We now see how SVM can be used for a restricted type of unsupervised learning, namely, for estimating regions of high density. We are not doing a full density estimation; rather, we want to find a boundary (so that it reads like a classification problem) that separates volumes of high density from volumes of low density (Tax and Duin 1999). Such a boundary can then be used for *novelty* or *outlier detection*. This is also called *one-class classification*.

OUTLIER DETECTION  
ONE-CLASS  
CLASSIFICATION

We consider a sphere with center  $\mathbf{a}$  and radius  $R$  that we want to enclose as much as possible of the density, measured empirically as the enclosed training set percentage. At the same time, trading off with it, we want to find the smallest radius (see figure 13.10). We define slack variables for instances that lie outside (we only have one type of slack variable because we have examples from one class and we do not have any penalty for those inside), and we have a smoothness measure that is proportional to the radius:

$$(13.52) \quad \min R^2 + C \sum_t \xi^t$$

subject to

$$\|\mathbf{x}^t - \mathbf{a}\|^2 \leq R^2 + \xi^t \text{ and } \xi^t \geq 0, \forall t$$

Adding the constraints, we get the Lagrangian, which we write keeping in mind that  $\|\mathbf{x}^t - \mathbf{a}\|^2 = (\mathbf{x}^t - \mathbf{a})^T (\mathbf{x}^t - \mathbf{a})$ :

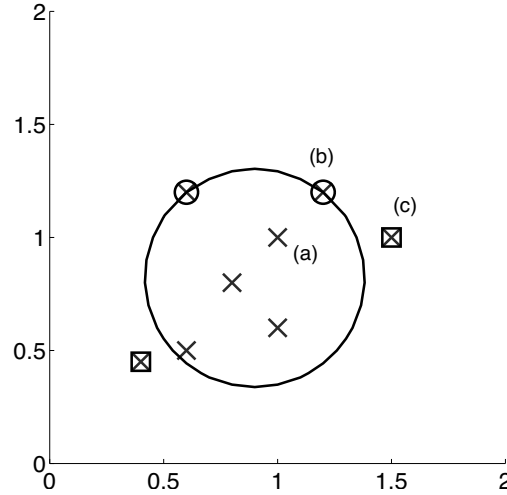
$$(13.53) \quad L_p = R^2 + C \sum_t \xi^t - \sum_t \alpha^t (R^2 + \xi^t - [(\mathbf{x}^t)^T \mathbf{x}^t - 2\mathbf{a}^T \mathbf{x}^t + \mathbf{a}^T \mathbf{a}]) - \sum_t \gamma^t \xi^t$$

with  $\alpha^t \geq 0$  and  $\gamma^t \geq 0$  being the Lagrange multipliers. Taking the derivative with respect to the parameters, we get

$$(13.54) \quad \frac{\partial L}{\partial R} = 2R - 2R \sum_t \alpha^t = 0 \Rightarrow \sum_t \alpha^t = 1$$

$$(13.55) \quad \frac{\partial L}{\partial \mathbf{a}} = \sum_t \alpha^t (2\mathbf{x}^t - 2\mathbf{a}) = 0 \Rightarrow \mathbf{a} = \sum_t \alpha^t \mathbf{x}^t$$

$$(13.56) \quad \frac{\partial L}{\partial \xi^t} = C - \alpha^t - \gamma^t = 0$$



**Figure 13.10** One-class support vector machine places the smoothest boundary (here using a linear kernel, the circle with the smallest radius) that encloses as much of the instances as possible. There are three possible cases: In (a), the instance is a typical instance. In (b), the instance falls on the boundary with  $\xi^t = 0$ ; such instances define  $R$ . In (c), the instance is an outlier with  $\xi^t > 0$ . (b) and (c) are support vectors. In terms of the dual variable, we have, in (a),  $\alpha^t = 0$ ; in (b),  $0 < \alpha^t < C$ ; in (c),  $\alpha^t = C$ .

Since  $\gamma^t \geq 0$ , we can write this last as the constraint:  $0 \leq \alpha^t \leq C$ . Plugging these into equation 13.53, we get the dual that we maximize with respect to  $\alpha^t$ :

$$(13.57) \quad L_d = \sum_t \alpha^t (\mathbf{x}^t)^T \mathbf{x}^t - \sum_t \sum_s \alpha^t \alpha^s (\mathbf{x}^t)^T \mathbf{x}^s$$

subject to

$$0 \leq \alpha^t \leq C \text{ and } \sum_t \alpha^t = 1$$

When we solve this, we again see that most of the instances vanish with their  $\alpha^t = 0$ ; these are the typical, highly likely instances that fall inside the sphere (figure 13.10). There are two type of support vectors with  $\alpha^t > 0$ : There are instances that satisfy  $0 < \alpha^t < C$  and lie on the boundary,  $\|\mathbf{x}^t - \mathbf{a}\|^2 = R^2$  ( $\xi^t = 0$ ), which we use to calculate  $R$ . Instances

that satisfy  $\alpha^t = C$  ( $\xi^t > 0$ ) lie outside the boundary and are the outliers. From equation 13.55, we see that the center  $\mathbf{a}$  is written as a weighted sum of the support vectors.

Then given a test input  $\mathbf{x}$ , we say that it is an outlier if

$$\|\mathbf{x} - \mathbf{a}\|^2 > R^2$$

or

$$\mathbf{x}^t \mathbf{x} - 2\mathbf{a}^T \mathbf{x} + \mathbf{a}^T \mathbf{a} > R^2$$

Using kernel functions, allow us to go beyond a sphere and define boundaries of arbitrary shapes. Replacing the dot product with a kernel function, we get (subject to the same constraints):

$$(13.58) \quad L_d = \sum_t \alpha^t K(\mathbf{x}^t, \mathbf{x}^t) - \sum_t \sum_s \alpha^t \alpha^s K(\mathbf{x}^t, \mathbf{x}^s)$$

For example, using a polynomial kernel of degree 2 allows arbitrary quadratic surfaces to be used. If we use a Gaussian kernel (equation 13.30), we have a union of local spheres. We reject  $\mathbf{x}$  as an outlier if

$$K(\mathbf{x}, \mathbf{x}) - 2 \sum_t \alpha^t K(\mathbf{x}, \mathbf{x}^t) + \sum_t \sum_s \alpha^t \alpha^s K(\mathbf{x}^t, \mathbf{x}^s) > R^2$$

The third term does not depend on  $\mathbf{x}$  and is therefore a constant (we use this as an equality to solve for  $R$  where  $\mathbf{x}$  is an instance on the margin). In the case of a Gaussian kernel where  $K(\mathbf{x}, \mathbf{x}) = 1$ , the condition reduces to

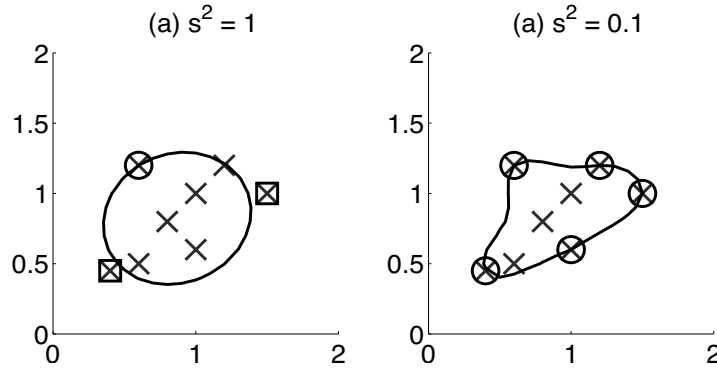
$$\sum_t \alpha^t K_G(\mathbf{x}, \mathbf{x}^t) < R_c$$

for some constant  $R_c$ , which is analogous to the kernel density estimator (section 8.2.2)—except for the sparseness of the solution—with a probability threshold  $R_c$  (see figure 13.11).

There is also an alternative, equivalent  $\nu$ -SVM type of formulation of one-class support vector machines that uses the canonical  $(1/2)\|\mathbf{w}\|^2$  type of smoothness (Schölkopf et al. 2001).

## 13.12 Kernel Dimensionality Reduction

We know from section 6.3 that principal components analysis (PCA) reduces dimensionality by projecting on the eigenvectors of the covariance



**Figure 13.11** One-class support vector machine using a Gaussian kernel with different spreads.

matrix  $\Sigma$  with the largest eigenvalues, which, if data instances are centered ( $E[\mathbf{x}] = 0$ ), can be written as  $\mathbf{X}^T \mathbf{X}$ . In the kernelized version, we work in the space of  $\phi(\mathbf{x})$  instead of the original  $\mathbf{x}$  and because, as usual, the dimensionality  $d$  of this new space may be much larger than the data set size  $N$ , we prefer to work with the  $N \times N$  matrix  $\mathbf{X}\mathbf{X}^T$  instead of the  $d \times d$  matrix  $\mathbf{X}^T \mathbf{X}$ . The projected data matrix is  $\Phi = \phi(\mathbf{X})$ , and hence we work on the eigenvectors of  $\Phi^T \Phi$  and hence of the kernel matrix  $\mathbf{K}$ .

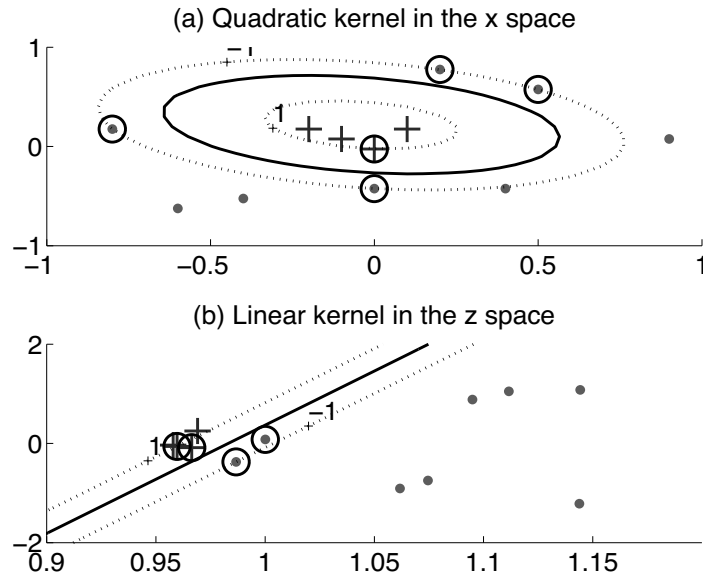
#### KERNEL PCA

*Kernel PCA* uses the eigenvectors and eigenvalues of the kernel matrix and this corresponds to doing a linear dimensionality reduction in the  $\phi(\mathbf{x})$  space. When  $c_i$  and  $\lambda_i$  are the corresponding eigenvectors and eigenvalues, the projected new  $k$ -dimensional values can be calculated as

$$\mathbf{z}_j^t = \sqrt{\lambda_i} c_j^t, j = 1, \dots, k, t = 1, \dots, N$$

An example is given in figure 13.12 where we first use a quadratic kernel and then decrease dimensionality to two (out of five) using kernel PCA and implement a linear SVM there. Note that in the general case (e.g., with a Gaussian kernel), the eigenvalues do not necessarily decay and there is no guarantee that we can reduce dimensionality using kernel PCA.

What we are doing here is multidimensional scaling (section 6.5) using kernel values as the similarity values. For example, by taking  $k = 2$ , one can visualize the data in the space induced by the kernel matrix, which can give us information as to how similarity is defined by the used kernel. Linear discriminality reduction (LDA) (section 6.6) can similarly



**Figure 13.12** Instead of using a quadratic kernel in the original space (a), we can use kernel PCA on the quadratic kernel values to map to a two-dimensional new space where we use a linear discriminant (b); these two dimensions (out of five) explain 80 percent of the variance.

be kernelized (Müller et al. 2001).

In chapter 6, we discussed nonlinear dimensionality reduction methods, Isomap and LLE. In fact, by viewing the elements of the cost matrix in equation 6.47 as kernel evaluations for pairs of inputs, LLE can be seen as kernel PCA for a particular choice of kernel. The same also holds for Isomap when a kernel function is defined as a function of the geodesic distance on the graph.

### 13.13 Notes

The idea of generalizing linear models by mapping the data to a new space through nonlinear basis functions is old, but the novelty of support vector machines is that of integrating this into a learning algorithm whose parameters are defined in terms of a subset of data instances (the so-called *dual representation*), hence also without needing to explicitly

evaluate the basis functions and thereby also limiting complexity by the size of the training set; this is also true for Gaussian processes where the kernel function is called the covariance function (section 14.4).

The sparsity of the solution shows the advantage over nonparametric estimators, such as  $k$ -nearest neighbor and Parzen windows, or Gaussian processes, and the flexibility to use kernel functions allows working with nonvectorial data. Because there is a unique solution to the optimization problem, we do not need any iterative optimization procedure as we do in neural networks. Because of all these reasons, support vector machines are now considered to be the best, off-the-shelf learners and are widely used in many domains, especially bioinformatics (Schölkopf, Tsuda, and Vert 2004) and natural language processing applications, where an increasing number of tricks are being developed to derive kernels (Shawe-Taylor and Cristianini 2004).

The use of kernel functions implies a different data representation; we no longer define an instance (object/event) as a vector of attributes by itself, but in terms of how it is similar to or differs from other instances; this is akin to the difference between multidimensional scaling that uses a matrix of distances (without any need to know how they are calculated) and principal components analysis that uses vectors in some space.

More information on support vector machines can be found in books by Vapnik (1995, 1998) and Schölkopf and Smola (2002). The chapter on SVM in Cherkassky and Mulier 1998 is a very readable introduction. Burges 1998 and Smola and Schölkopf 1998 are good tutorials on SVM classification and regression, respectively. There is a dedicated Web site <http://www.kernel-machines.org> and many free software packages are available, such as SVMlight (Joachims 2004) and LIBSVM (Chang and Lin 2008).

### 13.14 Exercises

1. Propose a filtering algorithm to find training instances that are very unlikely to be support vectors.
2. In equation 13.31, how can we estimate  $\mathbf{S}$ ?
3. In the empirical kernel map, how can we choose the templates?
4. In the localized multiple kernel of equation 13.40, propose a suitable model for  $\eta_i(\mathbf{x}|\theta_i)$  and discuss how it can be trained.
5. In kernel regression, what is the relation, if any, between  $\epsilon$  and noise variance?

6. In kernel regression, what is the effect of using different  $\epsilon$  on bias and variance?
7. How can we use one-class SVM for classification?
8. In a setting as in figure 13.12, use kernel PCA with a Gaussian kernel.
9. Let us say we have two representations for the same object and associated with each, we have a different kernel. How can we use both to implement a joint dimensionality reduction using kernel PCA?

## 13.15 References

- Allwein, E. L., R. E. Schapire, and Y. Singer. 2000. "Reducing Multiclass to Binary: A Unifying Approach for Margin Classifiers." *Journal of Machine Learning Research* 1: 113-141.
- Burges, C. J. C. 1998. "A Tutorial on Support Vector Machines for Pattern Recognition." *Data Mining and Knowledge Discovery* 2: 121-167.
- Chang, C.-C., and C.-J. Lin. 2008. *LIBSVM: A Library for Support Vector Machines*. <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.
- Cherkassky, V., and F. Mulier. 1998. *Learning from Data: Concepts, Theory, and Methods*. New York: Wiley.
- Cortes, C., and V. Vapnik. 1995. "Support Vector Networks." *Machine Learning* 20: 273-297.
- Dietterich, T. G., and G. Bakiri. 1995. "Solving Multiclass Learning Problems via Error-Correcting Output Codes." *Journal of Artificial Intelligence Research* 2: 263-286.
- Gönen, M., and E. Alpaydm. 2008. "Localized Multiple Kernel Learning." In *25th International Conference on Machine Learning*, ed. A. McCallum and S. Roweis, 352-359. Madison, WI: Omnipress.
- Jaakkola, T., and D. Haussler. 1999. "Exploiting Generative Models in Discriminative Classifiers." In *Advances in Neural Information Processing Systems 11*, ed. M. J. Kearns, S. A. Solla, and D. A. Cohn, 487-493. Cambridge, MA: MIT Press.
- Joachims, T. 2004. *SVMlight*, <http://svmlight.joachims.org>.
- Lanckriet, G. R. G., N. Cristianini, P. Bartlett, L. El Ghaoui, and M. I. Jordan. 2004. "Learning the Kernel Matrix with Semidefinite Programming." *Journal of Machine Learning Research* 5: 27-72.
- Mayoraz, E., and E. Alpaydm. 1999. "Support Vector Machines for Multiclass Classification." In *Foundations and Tools for Neural Modeling, Proceedings of IWANN'99, LNCS 1606*, ed. J. Mira and J. V. Sanchez-Andres, 833-842. Berlin: Springer.

- Müller, K. R., S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf. 2001. "An Introduction to Kernel-Based Learning Algorithms." *IEEE Transactions on Neural Networks* 12: 181–201.
- Noble, W. S. 2004. "Support Vector Machine Applications in Computational Biology." In *Kernel Methods in Computational Biology*, ed. B. Schölkopf, K. Tsuda, and J.-P. Vert, 71–92. Cambridge, MA: MIT Press.
- Platt, J. 1999. "Probabilities for Support Vector Machines." In *Advances in Large Margin Classifiers*, ed. A. J. Smola, P. Bartlett, B. Schölkopf, and D. Schuurmans, 61–74. Cambridge, MA: MIT Press.
- Schölkopf, B., J. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson. 2001. "Estimating the Support of a High-Dimensional Distribution." *Neural Computation* 13: 1443–1471.
- Schölkopf, B., and A. J. Smola. 2002. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA: MIT Press.
- Schölkopf, B., A. J. Smola, R. C. Williamson, and P. L. Bartlett. 2000. "New Support Vector Algorithms." *Neural Computation* 12: 1207–1245.
- Schölkopf, B., K. Tsuda, and J.-P. Vert, eds. 2004. *Kernel Methods in Computational Biology*. Cambridge, MA: MIT Press.
- Shawe-Taylor, J., and N. Cristianini. 2004. *Kernel Methods for Pattern Analysis*. Cambridge, UK: Cambridge University Press.
- Smola, A., and B. Schölkopf. 1998. *A Tutorial on Support Vector Regression*, NeuroCOLT TR-1998-030, Royal Holloway College, University of London, UK.
- Sonnenburg, S., G. Rätsch, C. Schäfer, and B. Schölkopf. 2006. "Large Scale Multiple Kernel Learning." *Journal of Machine Learning Research* 7: 1531–1565.
- Tax, D. M. J., and R. P. W. Duin. 1999. "Support Vector Domain Description." *Pattern Recognition Letters* 20: 1191–1199.
- Vapnik, V. 1995. *The Nature of Statistical Learning Theory*. New York: Springer.
- Vapnik, V. 1998. *Statistical Learning Theory*. New York: Wiley.
- Vert, J.-P., K. Tsuda, and B. Schölkopf. 2004. "A Primer on Kernel Methods." In *Kernel Methods in Computational Biology*, ed. B. Schölkopf, K. Tsuda, and J.-P. Vert, 35–70. Cambridge, MA: MIT Press.
- Weston, J., and C. Watkins. 1998. "Multiclass Support Vector Machines." *Technical Report CSD-TR-98-04*, Department of Computer Science, Royal Holloway, University of London.

# 14 *Bayesian Estimation*

*In the Bayesian approach, we consider parameters as random variables having a prior distribution. We continue from where we left off in section 4.4 and discuss three cases: estimating the parameters of a distribution, estimating the parameters of a model, and Gaussian processes.*

## 14.1 Introduction

BAYESIAN ESTIMATION is used when we have some prior information regarding a parameter. For example, before looking at a sample to estimate the mean  $\mu$  of a distribution, we may have some prior belief that it is close to 2, between 1 and 3. Such prior beliefs are especially important when we have a small sample. In such a case, we are interested in combining what the data tells us, namely, the value calculated from the sample, and our prior information.

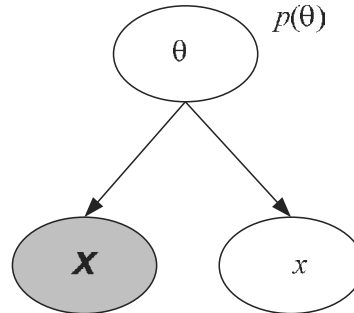
PRIOR PROBABILITY

The maximum likelihood approach we discuss in section 4.2 treats a parameter as an unknown constant. In Bayesian estimation, as we started discussing in section 4.4, a parameter is treated as a random variable, which allows us to code any prior information we have using a *prior probability distribution*. For example, knowing that  $\mu$  is very likely to be between 1 and 3, we write  $p(\mu)$  in such a way that the bulk of the density lies in the interval  $[1, 3]$ .

POSTERIOR  
PROBABILITY

Using Bayes' rule, we combine the prior and the likelihood and calculate the *posterior probability distribution*:

$$(14.1) \quad p(\theta|X) = \frac{p(\theta)p(X|\theta)}{p(X)}$$



**Figure 14.1** The generative graphical model. The arcs are in the direction of sampling; first we pick  $\theta$  from  $p(\theta)$  and then we generate data by sampling from  $p(x^t|\theta)$ . The new instance  $x$  and the past sample  $\mathcal{X}$  are independent given  $\theta$ : this is the iid assumption. If we do not know  $\theta$ , they are dependent: we infer  $\theta$  from given  $\mathcal{X}$  (shown shaded) using Bayes' rule, which inverts the direction to calculate  $p(\theta|\mathcal{X})$ , which can then be used to fill in  $x$ .

$p(\theta)$  is the prior density; it is what we know regarding the possible values that  $\theta$  may take *before* looking at the sample.  $p(\mathcal{X}|\theta)$  is the *sample likelihood*; it tells us how likely our sample  $\mathcal{X}$  is if the parameter of the distribution takes the value  $\theta$ . For example, if the instances in our sample are between 5 and 10, such a sample is likely if  $\mu$  is 7 but is less likely if  $\mu$  is 3 and even less likely if  $\mu$  is 1.  $p(\mathcal{X})$  in the denominator is a normalizer to make sure that the *posterior*  $p(\theta|\mathcal{X})$  integrates to 1. It is called the posterior probability because it tells us how likely  $\theta$  takes a certain value *after* looking at the sample. The Bayes' rule takes the prior distribution, combines it with what the data reveals, and generates the posterior distribution. We then use this posterior distribution in our later inferences.

For example, let us say that we have a past sample  $\mathcal{X}$  drawn from some distribution with unknown parameter  $\theta$ . We can then draw one more instance  $x$ , and we would like to calculate its probability distribution. We can visualize this as a graphical model (chapter 16) as shown in figure 14.1. What is depicted is a *generative model* which represents how the data is generated: We first pick  $\theta$  from  $p(\theta)$  and use it to sample  $\mathcal{X}$  and also the new instance  $x$ . We write the joint as

$$p(x, \mathcal{X}, \theta) = p(\theta)p(\mathcal{X}|\theta)p(x|\theta)$$

which we use in estimating the probability of a new instance  $x$  given the past sample  $\mathcal{X}$ :

$$\begin{aligned}
 p(x|\mathcal{X}) &= \frac{p(x, \mathcal{X})}{p(\mathcal{X})} = \frac{\int_{\theta} p(x, \mathcal{X}, \theta) d\theta}{p(\mathcal{X})} = \frac{\int_{\theta} p(\theta) p(\mathcal{X}|\theta) p(x|\theta) d\theta}{p(\mathcal{X})} \\
 (14.2) \qquad &= \int_{\theta} p(\theta|\mathcal{X}) p(x|\theta) d\theta
 \end{aligned}$$

In calculating  $p(\theta|\mathcal{X})$ , Bayes' rule allows us to invert the direction of the arc and do a diagnostic inference. The inferred  $\theta$  distribution is then used to derive a prediction distribution for the new  $x$ .

We see that our estimate is a weighted sum (we replace  $\int d\theta$  by  $\sum_{\theta}$  if  $\theta$  is discrete valued) of estimates for all possible values of  $\theta$  weighted by how likely  $\theta$  is, given the sample  $\mathcal{X}$ .

This is the *full Bayesian treatment* that may not be possible if the posterior is not easy to integrate. As we saw in section 4.4, in the case of the *maximum a posteriori (MAP) estimate*, we use the mode of the posterior:

MAXIMUM A  
POSTERIORI (MAP)  
ESTIMATE

$$\theta_{MAP} = \arg \max_{\theta} p(\theta|\mathcal{X}) \text{ and } p_{MAP}(x|\mathcal{X}) = p(x|\theta_{MAP})$$

The MAP estimate corresponds to assuming that the posterior makes a very narrow peak around a single point, that is, the mode. If the prior  $p(\theta)$  is uniform over all  $\theta$ , then the mode of the posterior  $p(\theta|\mathcal{X})$  and the mode of the likelihood  $p(\mathcal{X}|\theta)$  are at the same point, and the MAP estimate is equal to the maximum likelihood (ML) estimate. This implies that using ML corresponds to assuming no a priori distinction between different values of  $\theta$ .

Let us now see how Bayesian estimation is used in different types of distributions and applications.

## 14.2 Estimating the Parameter of a Distribution

### 14.2.1 Discrete Variables

Let us say that each instance is a multinomial variable taking one of  $K$  distinct states (section 4.2.2). We say  $x_i^t = 1$  if instance  $t$  is in state  $i$  and  $x_j^t = 0, \forall j \neq i$ . The parameters are the probabilities of states,  $\mathbf{q} = [q_1, q_2, \dots, q_k]^T$  with  $q_i, i = 1, \dots, K$  satisfying  $q_i \geq 0, \forall i$  and  $\sum_i q_i = 1$ .

The sample likelihood is

$$p(\mathcal{X}|\mathbf{q}) = \prod_{t=1}^N \prod_{i=1}^K q_i^{x_i^t}$$

DIRICHLET  
DISTRIBUTION

The prior distribution we use is the *Dirichlet distribution*

$$\text{Dirichlet}(\mathbf{q}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{i=1}^K q_i^{\alpha_i-1}$$

GAMMA FUNCTION

where  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_K]^T$  and  $\alpha_0 = \sum_i \alpha_i$ .  $\alpha_i$  being the parameters of the prior are called the *hyperparameters*.  $\Gamma(x)$  is the *Gamma function* defined as

$$\Gamma(x) \equiv \int_0^{\infty} u^{x-1} e^{-u} du$$

For example,  $x^t$  may correspond to news documents and states may correspond to  $K$  different news categories: sports, politics, arts, and so on. The probabilities  $q_i$  then correspond to the proportions of different news categories, and priors on them allow us to code our prior beliefs in these proportions; for example, we may expect to have more news related to sports than news related to arts.

Given the prior and the likelihood, we can derive the posterior

$$\begin{aligned} p(\mathbf{q}|\mathcal{X}) &\propto p(\mathcal{X}|\mathbf{q})p(\mathbf{q}|\boldsymbol{\alpha}) \\ (14.3) \quad &\propto \prod_i q_i^{\alpha_i+N_i-1} \end{aligned}$$

CONJUGATE PRIOR

where  $N_i = \sum_{t=1}^N x_i^t$ . We see that the posterior has the same form as the prior and we call such a prior a *conjugate prior*. Both the prior and the likelihood have the form of product of powers of  $q_i$ , and we combine them to make up the posterior:

$$\begin{aligned} p(\mathbf{q}|\mathcal{X}) &= \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + N_1) \cdots \Gamma(\alpha_K + N_K)} \prod_{i=1}^K q_i^{\alpha_i+N_i-1} \\ (14.4) \quad &= \text{Dirichlet}(\mathbf{q}|\boldsymbol{\alpha} + \mathbf{n}) \end{aligned}$$

where  $\mathbf{n} = [N_1, \dots, N_K]^T$  and  $\sum_i N_i = N$ .

Looking at equation 14.3, we can bring an interpretation to the hyperparameters  $\alpha_i$  (Bishop 2006). Just as  $n_i$  are counts of occurrences of state  $i$  in a sample of  $N$ , we can view  $\alpha_i$  as counts of occurrences of state  $i$  in some imaginary sample of  $\alpha_0$  instances. In defining the prior, we are subjectively saying the following: in a sample of  $\alpha_0$ , I would expect  $\alpha_i$  of them to belong to state  $i$ . Note that larger  $\alpha_0$  implies that we have a higher confidence (a more peaked distribution) in our subjective proportions: saying that I expect to have 60 out of 100 occurrences belong to

state 1 has higher confidence than saying that I expect to have 6 out of 10. The posterior then is another Dirichlet that sums up the counts of the occurrences of states, imagined and actual, given by the prior and the likelihood, respectively.

The conjugacy has a nice implication. In a sequential setting where we receive a sequence of instances, because the posterior and the prior have the same form, the current posterior accumulates information from all past instances and becomes the prior for the next instance.

When the variable is binary,  $x^t \in \{0, 1\}$ , the multinomial sample becomes Bernoulli

$$p(\mathcal{X}|q) = \prod_t q^{x^t} (1 - q)^{1-x^t}$$

BETA DISTRIBUTION

and the Dirichlet prior reduces to the *beta distribution*:

$$\text{beta}(q|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} q^{\alpha-1} (1 - q)^{\beta-1}$$

For example,  $x^t$  may be 0 or 1 depending on whether email with index  $t$  in a random sample of size  $N$  is legitimate or spam, respectively. Then defining a prior on  $q$  allows us to define a prior belief on the spam probability: I would expect, on the average,  $\alpha/(\alpha + \beta)$  of my emails to be spam.

Beta is a conjugate prior, and for the posterior we get

$$p(q|A, N, \alpha, \beta) \propto q^{A+\alpha-1} (1 - q)^{N-A+\beta-1}$$

where  $A = \sum_t x^t$ , and we see again that we combine the occurrences in the imaginary and the actual samples. Note that when  $\alpha = \beta = 1$ , we have a uniform prior and the posterior has the same shape as the likelihood. As the two counts, whether  $\alpha$  and  $\beta$  for the prior or  $\alpha + A$  and  $\beta + N - A$  for the posterior, increase and their difference increases, we get a distribution that is more peaked with smaller variance (see figure 14.2). As we see more data (imagined or actual), the variance decreases.

### 14.2.2 Continuous Variables

We now consider the case where instances are Gaussian distributed,  $p(x) \sim \mathcal{N}(\mu, \sigma^2)$ , and the parameters are  $\mu$  and  $\sigma^2$ ; we have already discussed this briefly section 4.4. The sample likelihood is

$$(14.5) \quad p(\mathcal{X}|\mu, \sigma^2) = \prod_t \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x^t - \mu)^2}{2\sigma^2}\right]$$

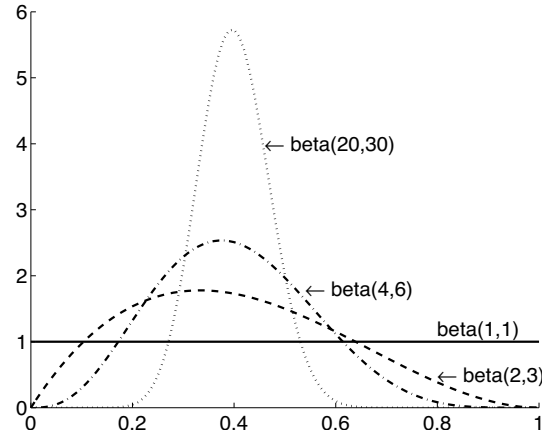


Figure 14.2 Plots of beta distributions for different sets of  $(\alpha, \beta)$ .

The conjugate prior for  $\mu$  is Gaussian,  $p(\mu) \sim \mathcal{N}(\mu_0^2, \sigma_0^2)$ , and we write the posterior as

$$\begin{aligned} p(\mu|\mathcal{X}) &\propto p(\mu)p(\mathcal{X}|\mu) \\ &\sim \mathcal{N}(\mu_N, \sigma_N^2) \end{aligned}$$

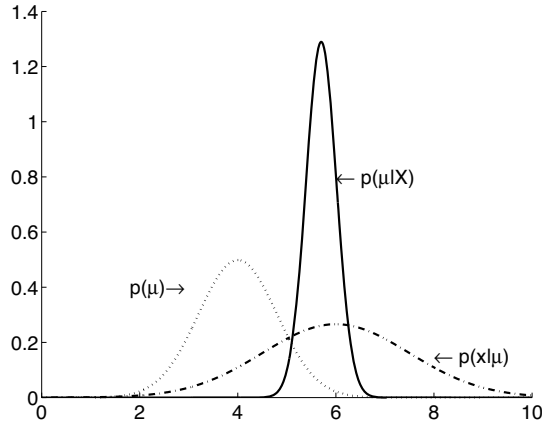
where

$$(14.6) \quad \mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}m$$

$$(14.7) \quad \frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

where  $m = \sum_t x^t / N$  is the sample average. We see that the mean of the posterior density (which is the Bayesian estimate),  $\mu_N$ , is a weighted average of the prior mean  $\mu_0$  and the sample mean  $m$ , with weights being inversely proportional to their variances (see figure 14.3 for an example). Note that because both coefficients are between 0 and 1 and sum to 1,  $\mu_N$  is always between  $\mu_0$  and  $m$ . When the sample size  $N$  or the variance of the prior  $\sigma_0^2$  is large, the Bayes' estimator is close to  $m$ , relying more on the information provided by the sample. When  $\sigma_0^2$  is small—that is, when we have little prior uncertainty regarding the correct value of  $\mu$ , or when we have a small sample—our prior guess  $\mu_0$  has higher effect.

$\sigma_N$  gets smaller when either of  $\sigma_0$  or  $\sigma$  gets smaller or if  $N$  is larger. Note also that  $\sigma_N$  is smaller than both  $\sigma_0$  and  $\sigma/\sqrt{N}$ , that is, the posterior



**Figure 14.3** 20 data points are drawn from  $p(x) \sim \mathcal{N}(6, 1.5^2)$ , prior is  $p(\mu) \sim \mathcal{N}(4, 0.8^2)$ , and posterior is then  $p(\mu|X) \sim \mathcal{N}(5.7, 0.3^2)$ .

variance is smaller than both prior variance and that of  $m$ . Incorporating both results in a better posterior estimate than using any of the prior or sample alone.

PRECISION

For the case of variance, we work with the *precision*, the reciprocal of the variance,  $\lambda \equiv 1/\sigma^2$ . Using this, the sample likelihood is written as

$$\begin{aligned}
 p(\mathcal{X}|\lambda) &= \prod_t \frac{\lambda^{1/2}}{\sqrt{2\pi}} \exp\left[-\frac{\lambda}{2}(x^t - \mu)^2\right] \\
 (14.8) \qquad &= \lambda^{N/2} (2\pi)^{-N/2} \exp\left[-\frac{\lambda}{2} \sum_t (x^t - \mu)^2\right]
 \end{aligned}$$

GAMMA DISTRIBUTION

The conjugate prior for the precision is the *Gamma distribution*:

$$p(\lambda) \sim \text{Gamma}(a_0, b_0) = \frac{1}{\Gamma(a_0)} b_0^{a_0} \lambda^{a_0-1} \exp(-b_0 \lambda)$$

and the posterior is

$$\begin{aligned}
 p(\lambda|\mathcal{X}) &\propto p(\mathcal{X}|\lambda)p(\lambda) \\
 &\sim \text{Gamma}(a_N, b_N)
 \end{aligned}$$

where

$$\begin{aligned}
 (14.9) \quad a_N &= a_0 + N/2 \\
 b_N &= b_0 + \frac{N}{2} s^2
 \end{aligned}$$

where  $s^2 = \sum_t (x^t - \mu)^2 / N$  is the sample variance. Again, we see that posterior estimates are weighted sum of priors and sample statistics.

### 14.3 Bayesian Estimation of the Parameters of a Function

We now discuss the case where we estimate the parameters, not of a distribution, but some function of the input, for regression or classification. Again, our approach is to consider these parameters as random variables with a prior distribution and use Bayes' rule to calculate a posterior distribution. We can then either evaluate the full integral, approximate it, or use the MAP estimate.

#### 14.3.1 Regression

Let us take the case of a linear regression model:

$$(14.10) \quad r = \mathbf{w}^T \mathbf{x} + \epsilon \text{ where } \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

where  $\beta$  is the precision of the additive noise.

The parameters are the weights  $\mathbf{w}$  and we have a sample  $\mathcal{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$  where  $\mathbf{x} \in \mathfrak{X}^d$  and  $r^t \in \mathfrak{R}$ , which we can break down into a matrix of inputs and a vector of desired outputs as  $\mathcal{X} = [\mathbf{X}, \mathbf{r}]$ . From equation 14.10, we have

$$p(r^t | \mathbf{x}^t, \mathbf{w}, \beta) \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}^t, \beta^{-1})$$

We saw previously in section 4.6 that the log likelihood is

$$\begin{aligned} \mathcal{L}(\mathcal{X} | \mathbf{w}) \equiv \log p(\mathcal{X} | \mathbf{w}) &= \log p(\mathbf{r}, \mathbf{X} | \mathbf{w}) \\ &= \log p(\mathbf{r} | \mathbf{X}, \mathbf{w}) + \log p(\mathbf{X}) \end{aligned}$$

where the second term is a constant, independent of the parameters. We expand the first term as

$$(14.11) \quad \begin{aligned} \mathcal{L}(\mathbf{r} | \mathbf{X}, \mathbf{w}, \beta) &= \log \prod_t p(r^t | \mathbf{x}^t, \mathbf{w}, \beta) \\ &= -N \log(\sqrt{2\pi}) + N \log \beta - \frac{\beta}{2} \sum_t (r^t - \mathbf{w}^T \mathbf{x}^t)^2 \end{aligned}$$

For the case of the ML estimate, we find  $\mathbf{w}$  that maximizes this, or equivalently, minimizes the last term that is the sum of the squared error, which can be rewritten as

$$E = (\mathbf{r} - \mathbf{X}\mathbf{w})^T (\mathbf{r} - \mathbf{X}\mathbf{w})$$

Taking the derivative with respect to  $\mathbf{w}$  and setting it to 0, we get the maximum likelihood estimator (we have previously derived this in section 5.8):

$$(14.12) \quad \mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r}$$

Having calculated the parameters, we can now do prediction. Given new input  $\mathbf{x}'$ , the response is calculated as

$$(14.13) \quad r' = \mathbf{w}_{ML}^T \mathbf{x}'$$

For nonlinear models,  $g(\mathbf{x}|\mathbf{w})$ , for example, a multilayer perceptron where  $\mathbf{w}$  are all the weights, we minimize, for example, using gradient descent,

$$E(\mathcal{X}|\mathbf{w}) = [r^t - g(\mathbf{x}^t|\mathbf{w})]^2$$

and  $\mathbf{w}_{LSQ}$  that minimize it is called the *least squares estimator*. Then, the prediction is calculated as

$$r' = g(\mathbf{x}'|\mathbf{w}_{LSQ})$$

GAUSSIAN PRIOR In the case of the Bayesian approach, for the parameters, we define a *Gaussian prior*:

$$p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$$

which is a conjugate prior and for the posterior, we get

$$p(\mathbf{w}|\mathcal{X}) \sim \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

where

$$(14.14) \quad \boldsymbol{\mu}_N = \beta \boldsymbol{\Sigma}_N \mathbf{X}^T \mathbf{r}$$

$$(14.15) \quad \boldsymbol{\Sigma}_N = (\alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X})^{-1}$$

To calculate the overall output, we integrate over the full posterior

$$r' = \int \mathbf{w}^T \mathbf{x}' p(\mathbf{w}|\mathcal{X}) d\mathbf{w}$$

If we want to use a point estimate, the MAP (or Bayes', because the posterior is Gaussian) estimator is

$$(14.16) \quad \mathbf{w}_{MAP} = \boldsymbol{\mu}_N = \beta (\alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{r}$$

and we replace the density with a single point, namely, the mean,

$$\mathbf{r}' = \mathbf{w}_{MAP}^T \mathbf{x}'$$

with variance

$$(14.17) \quad \text{Var}(\mathbf{r}') = \beta^{-1} + (\mathbf{x}')^T \Sigma_N \mathbf{x}'$$

Comparing equation 14.16 with the ML estimate of equation 14.12, this can be seen as regularization—that is, we add a constant  $\alpha$  to the diagonal to better condition the matrix to be inverted.

The prior,  $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I})$ , says that we expect the parameters to be close to 0 with spread inversely proportional to  $\alpha$ . When  $\alpha \rightarrow 0$ , we have a flat prior and the MAP estimate converges to the ML estimate.

We see in figure 14.4 that if we increase  $\alpha$ , we force parameters to be closer to 0 and the posterior distribution moves closer to the origin and shrinks. If we decrease  $\beta$ , we assume noise with higher variance and the posterior also has higher variance.

If we take the log of the posterior, we have

$$\begin{aligned} \log p(\mathbf{w}|\mathbf{X}, \mathbf{r}) &\propto \log p(\mathbf{X}, \mathbf{r}|\mathbf{w}) + \log p(\mathbf{w}) \\ &\propto \log p(\mathbf{r}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}) \\ &= -\frac{\beta}{2} \sum_t (r^t - \mathbf{w}^T \mathbf{x}^t)^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + c \end{aligned}$$

which we maximize to find the MAP estimate. In the general case, given our model  $g(\mathbf{x}|\mathbf{w})$ , we can write an augmented error function

$$E_{\text{ridge}}(\mathbf{w}|\mathcal{X}) = \sum_t [r^t - g(\mathbf{x}^t|\mathbf{w})]^2 + \lambda \sum_i w_i^2$$

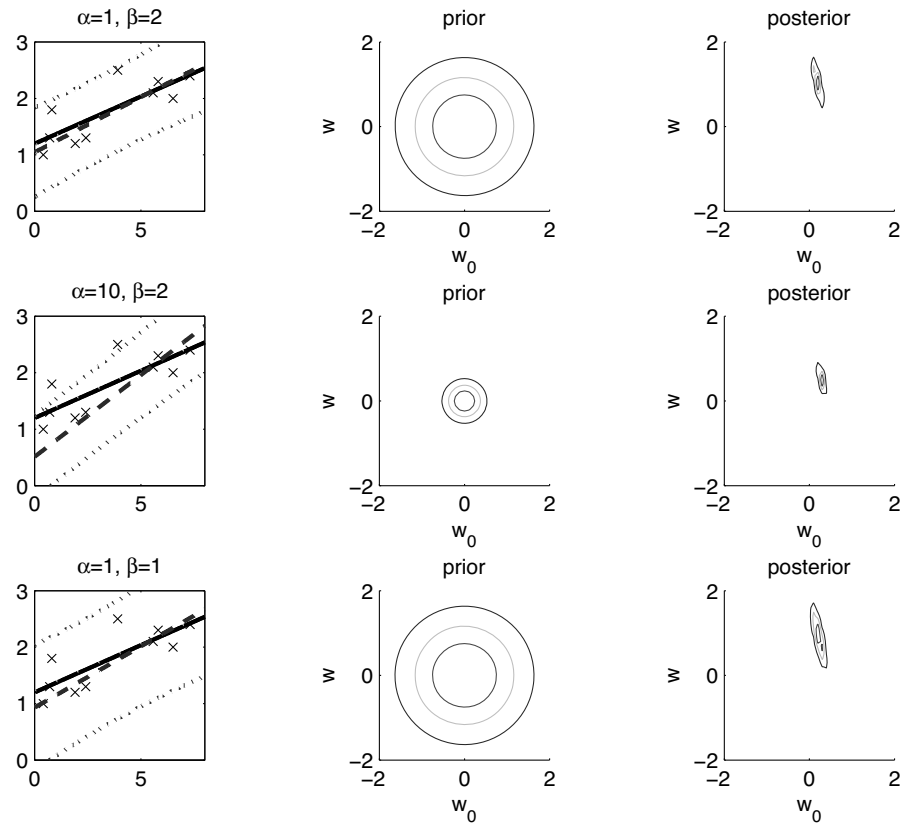
RIDGE REGRESSION

with  $\lambda \equiv \alpha/\beta$ . This is known as *parameter shrinkage* or *ridge regression* in statistics. In section 4.8, we called this *regularization* and in section 11.9, we called this *weight decay* in neural networks. The first term is the negative log of the likelihood, and the second term penalizes  $w_i$  away from 0 (as dictated by  $\alpha$  of the prior).

LAPLACIAN PRIOR

Though this approach reduces  $\sum_i w_i^2$ , it does not force individual  $w_i$  to 0; that is, it cannot be used for feature selection, namely, to determine which  $x_i$  are redundant. For this, one can use a *Laplacian prior* that uses the  $L_1$  norm instead of the  $L_2$  norm (Figueiredo 2003):

$$p(\mathbf{w}|\alpha) = \prod_i \frac{\alpha}{2} \exp(-\alpha|w_i|) = \left(\frac{\alpha}{2}\right)^d \exp\left(-\alpha \sum_i |w_i|\right)$$



**Figure 14.4** Bayesian linear regression for different values of  $\alpha$  and  $\beta$ . To the left: crosses are the data points and straight line is the ML solution. The MAP solution with one standard deviation error bars are also shown dashed. Center: prior density centered at 0 and variance  $1/\alpha$ . To the right: posterior density whose mean is the MAP solution. We see that when  $\alpha$  is increased, the variance of the prior shrinks and the line moves closer to the flat 0 line. When  $\beta$  is decreased, more noise is assumed and the posterior density has higher variance.

The posterior probability is no longer Gaussian and the MAP estimate is found by minimizing

$$E_{\text{lasso}}(\mathbf{w}|\mathcal{X}) = \sum_t (r^t - \mathbf{w}^T \mathbf{x}^t)^2 + 2\sigma^2 \alpha \sum_i |w_i|$$

where  $\sigma^2$  is the variance of noise (for which we plug in our estimate).

LASSO

This is known as *lasso* (least absolute shrinkage and selection operator) (Tibshirani 1996). To see why  $L_1$  induces sparseness, let us consider the case with two weights  $[w_1, w_2]^T$  (Figueiredo 2003):  $\|[1, 0]^T\|_2 = \|[1/\sqrt{2}, 1/\sqrt{2}]^T\|_2 = 1$ , whereas  $\|[1, 0]^T\|_1 = 1 < \|[1/\sqrt{2}, 1/\sqrt{2}]^T\|_1 = \sqrt{2}$ , and therefore  $L_1$  prefers to set  $w_2$  to 0 and use a large  $w_1$ , rather than having small values for both.

### 14.3.2 The Use of Basis/Kernel Functions

Using the Bayes' estimate of equation 14.14, the prediction is written as

$$\begin{aligned} r' &= (\mathbf{x}')^T \mathbf{w} \\ &= \beta (\mathbf{x}')^T \Sigma_N \mathbf{X}^T \mathbf{r} \\ &= \sum_t \beta (\mathbf{x}')^T \Sigma_N \mathbf{x}^t r^t \end{aligned}$$

DUAL REPRESENTATION

This is the *dual representation*. When we can write the parameter in terms of the training data, or a subset of it as in support vector machines (chapter 13), we can write the prediction as a function of the current input and past data. We can rewrite this as

$$(14.18) \quad r' = \sum_t K(\mathbf{x}', \mathbf{x}^t) r^t$$

where we define

$$(14.19) \quad K(\mathbf{x}', \mathbf{x}^t) = \beta (\mathbf{x}')^T \Sigma_N \mathbf{x}^t$$

BASIS FUNCTION

We know that we can generalize the linear kernel of equation 14.19 by using a nonlinear *basis function*  $\phi(\mathbf{x})$  to map to a new space where we fit the linear model. In such a case, instead of the  $d$ -dimensional  $\mathbf{x}$  we have the  $k$ -dimensional  $\phi(\mathbf{x})$  where  $k$  is the number of basis functions and instead of  $N \times d$  data matrix  $\mathbf{X}$ , we have  $N \times k$  image of the basis functions  $\Phi$ .

During test, we have

$$r' = \phi(\mathbf{x}')^T \mathbf{w} \text{ where } \mathbf{w} = \beta \Sigma_N^\phi \Phi^T \mathbf{r} \text{ and } \Sigma_N^\phi = (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1}$$

$$\begin{aligned}
 &= \beta \boldsymbol{\phi}(\mathbf{x}')^T \boldsymbol{\Sigma}_N^\phi \boldsymbol{\Phi}^T \mathbf{r} \\
 &= \sum_t \beta \boldsymbol{\phi}(\mathbf{x}')^T \boldsymbol{\Sigma}_N^\phi \boldsymbol{\phi}(\mathbf{x}^t) r^t \\
 (14.20) \quad &= \sum_t K(\mathbf{x}', \mathbf{x}^t) r^t
 \end{aligned}$$

where we define

$$(14.21) \quad K(\mathbf{x}', \mathbf{x}^t) = \beta \boldsymbol{\phi}(\mathbf{x}')^T \boldsymbol{\Sigma}_N^\phi \boldsymbol{\phi}(\mathbf{x}^t)$$

KERNEL FUNCTION

as the equivalent kernel. This is the dual representation in the space of  $\boldsymbol{\phi}(\mathbf{x})$ . We see that we can write our estimate as a weighted sum of the effects of instances in the training set where the effect is given by the *kernel function*  $K(\mathbf{x}', \mathbf{x}^t)$ ; this is similar to nonparametric kernel smoothers we discuss in chapter 8, or kernel machines of chapter 13.

Error bars can be defined using

$$\text{Var}(r') = \beta^{-1} + \boldsymbol{\phi}(\mathbf{x}')^T \boldsymbol{\Sigma}_N^\phi \boldsymbol{\phi}(\mathbf{x}')$$

An example is given in figure 14.5 for the linear, quadratic, and fourth-degree kernels.

Just as in regression proper where we can work on the original  $\mathbf{x}$  or  $\boldsymbol{\phi}(\mathbf{x})$ , in Bayesian regression too we can work on the preprocessed  $\boldsymbol{\phi}(\mathbf{x})$ , defining parameters in that space. Later on in this chapter, we are going to see Gaussian processes where we can define and use  $K(\mathbf{x}, \mathbf{x}^t)$  directly without needing to calculate  $\boldsymbol{\phi}(\mathbf{x})$ .

### 14.3.3 Bayesian Classification

In a two-class problem, we have a single output and assuming a linear model, we have

$$P(C_1 | \mathbf{x}^t) = y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t)$$

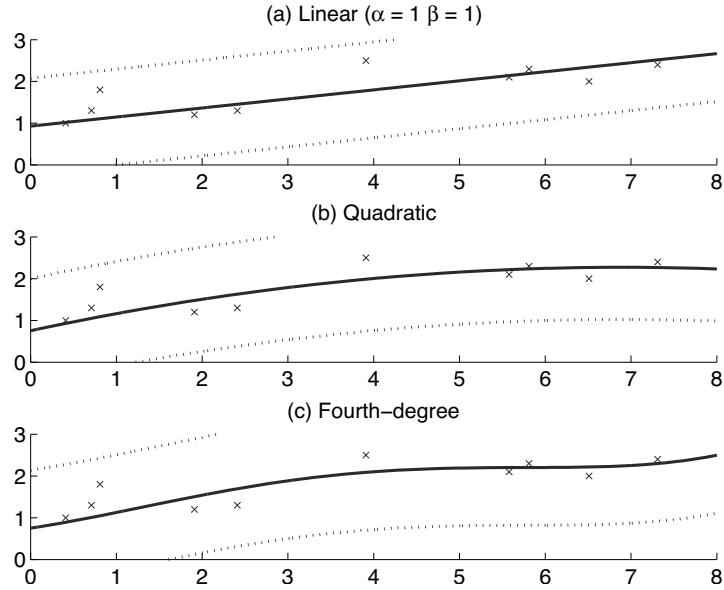
The log likelihood of a Bernoulli sample is given as

$$\mathcal{L}(\mathbf{r} | \mathbf{X}) = \sum_t r^t \log y_t + (1 - r^t) \log(1 - y^t)$$

which we maximize, or minimize its negative log—the cross-entropy—to find the ML estimate, for example, using gradient descent. This is called *logistic discrimination* (section 10.7).

In the case of the Bayesian approach, we assume a Gaussian prior

$$(14.22) \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$$



**Figure 14.5** Bayesian regression using kernels with one standard deviation error bars: (a) linear:  $\phi(x) = [1, x]^T$ , (b) quadratic:  $\phi(x) = [1, x, x^2]^T$ , and (c) fourth degree:  $\phi(x) = [1, x, x^2, x^3, x^4]^T$ .

and the log of the posterior is given as

$$\begin{aligned}
 \log p(\mathbf{w}|\mathbf{r}, \mathbf{X}) &\propto \log p(\mathbf{w}) + \log p(\mathbf{r}|\mathbf{w}, \mathbf{X}) \\
 &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \\
 &\quad + \sum_t r^t \log y_t + (1 - r^t) \log(1 - y^t) + c
 \end{aligned}
 \tag{14.23}$$

LAPLACE APPROXIMATION

This posterior distribution is no longer Gaussian and we cannot integrate exactly. We can use *Laplace approximation*, which works as follows (MacKay 2003). Let us say we want to approximate some distribution  $f(x)$ , not necessarily normalized (to integrate to 1). In Laplace approximation, we find the mode of  $f(x)$ ,  $x_0$ , fit a Gaussian  $q(x)$  centered there, and then if we want to integrate, we integrate this fitted Gaussian instead. To find the variance of the Gaussian, we consider the Taylor expansion of  $f(\cdot)$  at  $x = x_0$

$$\log f(x) = \log f(x_0) - \frac{1}{2}a(x - x_0)^2 + \dots$$

where

$$a \equiv - \left. \frac{d}{dx^2} \log f(x) \right|_{x=x_0}$$

Note that the first, linear term disappears because the first derivative is 0 at the mode. Taking exp, we have

$$f(x) = f(x_0) \exp \left[ -\frac{a}{2}(x - x_0)^2 \right]$$

To normalize  $f(x)$ , we consider that in a Gaussian distribution

$$\int \frac{1}{\sqrt{2\pi}(1/\sqrt{a})} \exp \left[ -\frac{a}{2}(x - x_0)^2 \right] = 1 \Rightarrow \int \exp \left[ -\frac{a}{2}(x - x_0)^2 \right] = \sqrt{a/2\pi}$$

and therefore

$$q(x) = \sqrt{a/2\pi} \exp \left[ -\frac{a}{2}(x - x_0)^2 \right] \sim \mathcal{N}(x_0, 1/a)$$

In the multivariate setting where  $\mathbf{x} \in \mathfrak{R}^d$ , we have

$$\log f(\mathbf{x}) = \log f(\mathbf{x}_0) - \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{A}(\mathbf{x} - \mathbf{x}_0) + \dots$$

where  $\mathbf{A}$  is the (Hessian) matrix of second derivatives:

$$\mathbf{A} = - \nabla \nabla \log f(\mathbf{x}) \big|_{\mathbf{x}=\mathbf{x}_0}$$

The Laplace approximation is then

$$f(\mathbf{x}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{d/2}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{A}(\mathbf{x} - \mathbf{x}_0) \right] \sim \mathcal{N}_d(\mathbf{x}_0, \mathbf{A}^{-1})$$

Having now discussed how to approximate, we can now use it for the posterior density.  $\mathbf{w}_{MAP}$ , which is the mode of  $p(\mathbf{w}|\mathbf{r}, \mathbf{X})$ , is taken as the mean and the covariance matrix is given by the inverse of the matrix of the second derivatives of the negative log likelihood:

$$\mathbf{S}_N = -\nabla \nabla \log p(\mathbf{w}|\mathbf{r}, \mathbf{X}) = \mathbf{S}_0^{-1} + \sum_t y^t (1 - y^t) \mathbf{x}^t (\mathbf{x}^t)^T$$

We then integrate over this Gaussian to estimate the class probability:

$$P(C_1|\mathbf{x}) = y = \int \text{sigmoid}(\mathbf{w}^T \mathbf{x}) q(\mathbf{w}) d\mathbf{w}$$

where  $q(\mathbf{w}) \sim \mathcal{N}(\mathbf{w}_{MAP}, \mathbf{S}_N^{-1})$ . A further complication is that we cannot integrate analytically over a Gaussian convolved with a sigmoid. If we use the *probit function* instead, which has the same S-shape as the sigmoid, an analytical solution is possible (Bishop 2006).

## 14.4 Gaussian Processes

Let us say we have the linear model  $y = \mathbf{w}^T \mathbf{x}$ . Then, for each  $\mathbf{w}$ , we have one line. Given a prior distribution  $p(\mathbf{w})$ , we get a distribution of lines, or to be more specific, for any  $\mathbf{w}$ , we get a distribution of  $y$  values calculated at  $\mathbf{x}$  as  $y(\mathbf{x}|\mathbf{w})$  when  $\mathbf{w}$  is sampled from  $p(\mathbf{w})$ , and this is what we mean when we talk about a Gaussian process. We know that if  $p(\mathbf{w})$  is Gaussian, each  $y$  is a linear combination of Gaussians and is also Gaussian; in particular, we are interested in the joint distribution of  $y$  values calculated at the  $N$  input data points,  $\mathbf{x}^t, t = 1, \dots, N$  (MacKay 1998).

We assume a zero mean Gaussian prior

$$p(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I})$$

Given the  $N \times d$  data points  $\mathbf{X}$  and the  $d \times 1$  weight vector, we write the  $y$  outputs as

$$(14.24) \quad \mathbf{y} = \mathbf{X}\mathbf{w}$$

which is  $N$ -variate Gaussian with

$$(14.25) \quad \begin{aligned} E[\mathbf{y}] &= \mathbf{X}E[\mathbf{w}] = \mathbf{0} \\ \text{Cov}(\mathbf{y}) &= E[\mathbf{y}\mathbf{y}^T] = \mathbf{X}E[\mathbf{w}\mathbf{w}^T]\mathbf{X}^T = \frac{1}{\alpha}\mathbf{X}\mathbf{X}^T \equiv \mathbf{K} \end{aligned}$$

where  $\mathbf{K}$  is the (Gram) matrix with elements

$$K_{i,j} \equiv K(\mathbf{x}^i, \mathbf{x}^j) = \frac{(\mathbf{x}^i)^T \mathbf{x}^j}{\alpha}$$

COVARIANCE  
FUNCTION

This is known as the *covariance function* in the literature of Gaussian processes and the idea is the same as in kernel functions: If we use a set of basis functions  $\boldsymbol{\phi}(\mathbf{x})$ , we generalize from the dot product of the original inputs to the dot product of basis functions by a kernel

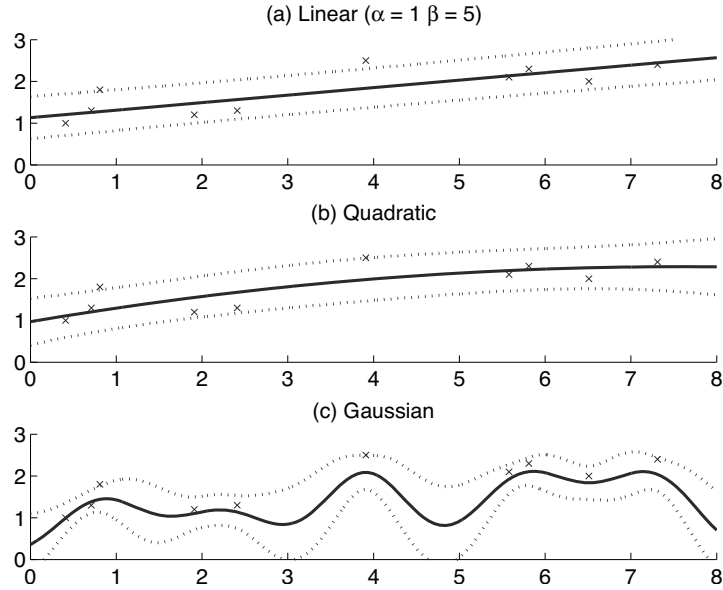
$$K_{i,j} = \frac{\boldsymbol{\phi}(\mathbf{x}^i)^T \boldsymbol{\phi}(\mathbf{x}^j)}{\alpha}$$

The actual observed output  $r$  is given by the line with added noise,  $r = y + \epsilon$  where  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ . For all  $N$  data points, we write it as

$$(14.26) \quad \mathbf{r} \sim \mathcal{N}_N(\mathbf{0}, \mathbf{C}_N) \text{ where } \mathbf{C}_N = \beta^{-1}\mathbf{I} + \mathbf{K}$$

To make a prediction, we consider the new data as the  $(N + 1)$ st data point pair  $(\mathbf{x}', r')$ , and write the joint using all  $N + 1$  data points. We have

$$(14.27) \quad \mathbf{r}_{N+1} \sim \mathcal{N}_N(\mathbf{0}, \mathbf{C}_{N+1})$$



**Figure 14.6** Gaussian process regression with one standard deviation error bars: (a) linear kernel, (b) quadratic kernel, (c) Gaussian kernel with spread  $s^2 = 0.5$ .

where

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{bmatrix}$$

with  $\mathbf{k}$  being the  $N \times 1$  dimensional vector of  $K(\mathbf{x}', \mathbf{x}^t), t = 1, \dots, N$  and  $c = K(\mathbf{x}', \mathbf{x}') + \beta^{-1}$ . Then to make a prediction, we calculate  $p(r' | \mathbf{x}', \mathbf{X}, \mathbf{r})$ , which is Gaussian with

$$\begin{aligned} E[r' | \mathbf{x}'] &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{r} \\ \text{Var}(r' | \mathbf{x}') &= c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \end{aligned}$$

An example is shown in figure 14.6 using linear, quadratic, and Gaussian kernels. The first two are defined as the dot product of their corresponding basis functions; the Gaussian kernel is defined directly as

$$K_G(\mathbf{x}^i, \mathbf{x}^j) = \exp \left[ -\frac{\|\mathbf{x}^i - \mathbf{x}^j\|^2}{s^2} \right]$$

The mean, which is our point estimate (if we do not integrate over the full distribution), can also be written as a weighted sum of the kernel effects

$$(14.28) \quad E[r' | \mathbf{x}'] = \sum_t a^t K(\mathbf{x}^t, \mathbf{x}')$$

where  $a^t$  is the  $t$ th component of  $\mathbf{C}_N^{-1} \mathbf{r}$ . Or, we can write it as a weighted sum of the outputs of the training data points where weights are given by the kernel function

$$(14.29) \quad E[r' | \mathbf{x}'] = \sum_t r^t w^t$$

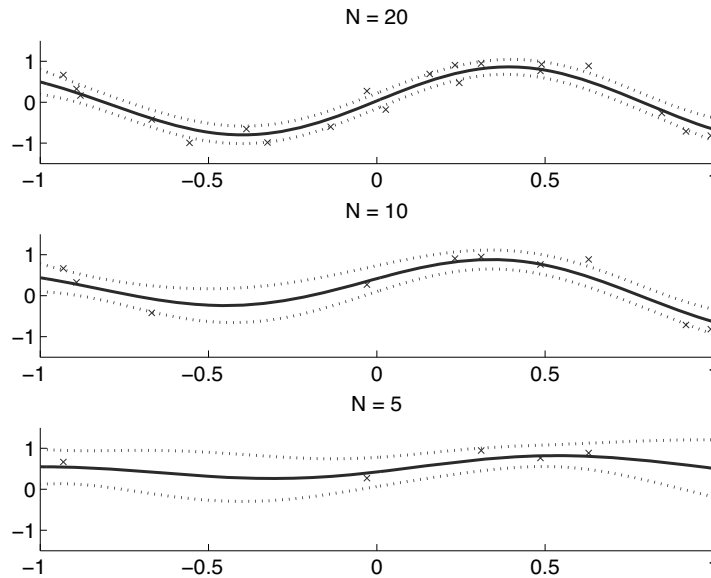
where  $w^t$  is the  $t$ th component of  $\mathbf{k}^T \mathbf{C}_N^{-1}$ .

Note that we can also calculate the variance of a prediction at a point to get an idea about uncertainty in there, and it depends on the instances that affect the prediction in there. In the case of a Gaussian kernel, only instances within a locality are effective and prediction variance is high where there is little data in the vicinity (see figure 14.7).

Kernel functions can be defined and used, depending on the application, as we have previously discussed in the context of kernel machines in chapter 13. The possibility of using kernel functions directly without needing to calculate or store the basis functions offers a great flexibility. Normally, given a training set, we first calculate the parameters, for example using equation 14.12, and then use the parameters to make predictions using equation 14.13, never needing the training set any more. This makes sense because generally the dimensionality of the parameters, which is generally  $\mathcal{O}(d)$ , is much lower than the size of the training set  $N$ .

When we work with basis functions, however, calculating the parameter explicitly may no longer be the case, because the dimensionality of the basis functions may be very high, even infinite. In such a case, it is cheaper to use the dual representation, taking into account the effects of training instances using kernel functions, as we do here. This idea is also used in nonparametric smoothers (chapter 8) and kernel machines (chapter 13).

The requirement here is that  $\mathbf{C}_N$  be invertible and hence positive definite. For this,  $\mathbf{K}$  should be semidefinite so that after adding  $\beta^{-1} > 0$  to the diagonals, we get positive definiteness. We also see that the costliest operation is this inversion of  $N \times N$  matrix, which fortunately needs to be



**Figure 14.7** Gaussian process regression using a Gaussian kernel with  $s^2 = 0.5$  and varying number of training data. We see how variance of the prediction is larger where there is few data.

calculated only once (during training) and stored. Still, for large  $N$ , one may need an approximation.

When we use it for classification for a two-class problem, the output is filtered through a sigmoid,  $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x})$ , and the distribution of  $y$  is no longer Gaussian. The derivation is similar except that the conditional  $p(r_{N+1} | \mathbf{x}_{N+1}, \mathbf{X}, \mathbf{r})$  is not Gaussian either and we need to approximate, for example, using Laplace approximation (Bishop 2006; Rasmussen and Williams 2006).

## 14.5 Notes

Bayesian approaches have become popular recently with advances in computational power allowing us to sample from or approximate the posterior probabilities. Truth has many cloaks. This preference of simplicity appears in many contexts as the Bayesian approach, regularization, min-

imum description length, or smoothing, and is at the heart of statistical inference and hence machine learning.

On the other hand, the subjectivity of priors is disturbing and there are objections to the Bayesian approach; see Gelman 2008, for example. What is the use of a flat prior, and why collect data if we already have a peaked prior? Is a conjugate prior true or merely convenient?

Just like with support vector machines, in Gaussian processes too, there are methods by which one can construct new kernels as functions (e.g., weighted sums) of some other kernels and these weights or kernel parameters (e.g., spreads) can be optimized by a *type 2 maximum likelihood procedure*, so called because we are now optimizing not the parameters (which are the  $a^t$  or  $w^t$  above) but the hyperparameters on a second level (Bishop 2006; Rasmussen and Williams 2006).

TYPE 2 MAXIMUM  
LIKELIHOOD  
PROCEDURE

## 14.6 Exercises

1. For the setting of figure 14.3, observe how the posterior changes as we change  $N$ ,  $\sigma^2$ , and  $\sigma_0^2$ .
2. Let us denote by  $x$  the number of spam emails I receive in a random sample of  $n$ . Assume that the prior for  $q$ , the proportion of spam emails is uniform in  $[0, 1]$ . Find the posterior distribution for  $p(q|x)$ .
3. As above, except that assume that  $p(q) \sim \mathcal{N}(\mu_0, \sigma_0^2)$ . Also assume  $n$  is large so that you can use central limit theorem and approximate binomial by a Gaussian. Derive  $p(q|x)$ .
4. What is  $\text{Var}(r')$  when the maximum likelihood estimator is used? Compare it with equation 14.17.
5. In figure 14.6, how does the fit change when we change  $s^2$ ?
6. Propose a filtering algorithm to choose a subset of the training set in Gaussian processes.
7. *Active learning* is when the learner is able to generate  $x$  itself and ask a supervisor to provide the corresponding  $r$  value during learning one by one, instead of passively being given a training set. How can we implement active learning using Gaussian processes? (Hint: Where do we have the largest uncertainty?)
8. Let us say we have inputs from two different representations. How can we use the approaches discussed in this chapter in such a case?

ACTIVE LEARNING

## 14.7 References

- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning*. New York: Springer.
- Figueiredo, M. A. T. 2003. "Adaptive Sparseness for Supervised Learning." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 25: 1150–1159.
- Gelman, A. 2008. "Objections to Bayesian statistics." *Bayesian Statistics* 3: 445–450.
- MacKay, D. J. C. 1998. "Introduction to Gaussian Processes." In *Neural Networks and Machine Learning*, ed. C. M. Bishop, 133–166. Berlin: Springer.
- MacKay, D. J. C. 2003. *Information Theory, Inference, and Learning Algorithms*. Cambridge, UK: Cambridge University Press.
- Rasmussen, C. E. , and C. K. I. Williams. 2006. *Gaussian Processes for Machine Learning*. Cambridge, MA: MIT Press.
- Tibshirani, R. 1996. "Regression Shrinkage and Selection via the Lasso." *Journal of the Royal Statistical Society B* 58: 267–288.



# 15 *Hidden Markov Models*

*We relax the assumption that instances in a sample are independent and introduce Markov models to model input sequences as generated by a parametric random process. We discuss how this modeling is done as well as introduce an algorithm for learning the parameters of such a model from example sequences.*

## 15.1 Introduction

UNTIL NOW, we assumed that the instances that constitute a sample are iid. This has the advantage that the likelihood of the sample is simply the product of the likelihoods of the individual instances. This assumption, however, is not valid in applications where successive instances are dependent. For example, in a word successive letters are dependent; in English ‘h’ is very likely to follow ‘t’ but not ‘x’. Such processes where there is a *sequence* of observations—for example, letters in a word, base pairs in a DNA sequence—cannot be modeled as simple probability distributions. A similar example is speech recognition where speech utterances are composed of speech primitives called phonemes; only certain sequences of phonemes are allowed, which are the words of the language. At a higher level, words can be written or spoken in certain sequences to form a sentence as defined by the syntactic and semantic rules of the language.

A sequence can be characterized as being generated by a *parametric random process*. In this chapter, we discuss how this modeling is done and also how the parameters of such a model can be learned from a training sample of example sequences.

## 15.2 Discrete Markov Processes

Consider a system that at any time is in one of a set of  $N$  distinct states:  $S_1, S_2, \dots, S_N$ . The state at time  $t$  is denoted as  $q_t, t = 1, 2, \dots$ , so, for example,  $q_t = S_i$  means that at time  $t$ , the system is in state  $S_i$ . Though we write “time” as if this should be a temporal sequence, the methodology is valid for any sequencing, be it in time, space, position on the DNA string, and so forth.

At regularly spaced discrete times, the system moves to a state with a given probability, depending on the values of the previous states:

$$P(q_{t+1} = S_j | q_t = S_i, q_{t-1} = S_k, \dots)$$

MARKOV MODEL

For the special case of a first-order *Markov model*, the state at time  $t + 1$  depends only on state at time  $t$ , regardless of the states in the previous times:

$$(15.1) \quad P(q_{t+1} = S_j | q_t = S_i, q_{t-1} = S_k, \dots) = P(q_{t+1} = S_j | q_t = S_i)$$

This corresponds to saying that, given the present state, the future is independent of the past. This is just a mathematical version of the saying, Today is the first day of the rest of your life.

We further simplify the model—that is, regularize—by assuming that these *transition probabilities* are independent of time:

TRANSITION  
PROBABILITIES

$$(15.2) \quad a_{ij} \equiv P(q_{t+1} = S_j | q_t = S_i)$$

satisfying

$$(15.3) \quad a_{ij} \geq 0 \text{ and } \sum_{j=1}^N a_{ij} = 1$$

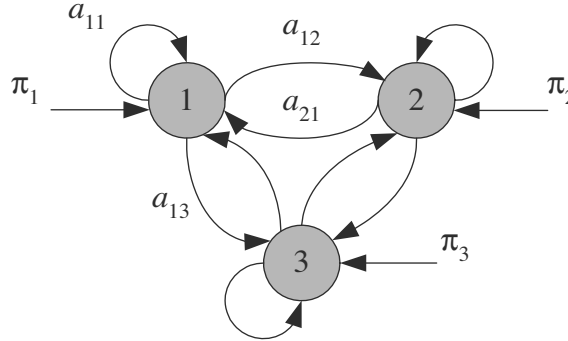
So, going from  $S_i$  to  $S_j$  has the same probability no matter when it happens, or where it happens in the observation sequence.  $\mathbf{A} = [a_{ij}]$  is a  $N \times N$  matrix whose rows sum to 1.

STOCHASTIC  
AUTOMATON

This can be seen as a *stochastic automaton* (see figure 15.1). From each state  $S_i$ , the system moves to state  $S_j$  with probability  $a_{ij}$ , and this probability is the same for any  $t$ . The only special case is the first state. We define *initial probabilities*,  $\pi_i$ , which is the probability that the first state in the sequence is  $S_i$ :

INITIAL PROBABILITIES

$$(15.4) \quad \pi_i \equiv P(q_1 = S_i)$$



**Figure 15.1** Example of a Markov model with three states. This is a stochastic automaton where  $\pi_i$  is the probability that the system starts in state  $S_i$ , and  $a_{ij}$  is the probability that the system moves from state  $S_i$  to state  $S_j$ .

satisfying

$$(15.5) \quad \sum_{i=1}^N \pi_i = 1$$

$\Pi = [\pi_i]$  is a vector of  $N$  elements that sum to 1.

OBSERVABLE MARKOV MODEL

In an *observable Markov model*, the states are observable. At any time  $t$ , we know  $q_t$ , and as the system moves from one state to another, we get an observation sequence that is a sequence of states. The output of the process is the set of states at each instant of time where each state corresponds to a physical observable event.

We have an observation sequence  $O$  that is the state sequence  $O = Q = \{q_1 q_2 \cdots q_T\}$ , whose probability is given as

$$(15.6) \quad P(O = Q | \mathbf{A}, \Pi) = P(q_1) \prod_{t=2}^T P(q_t | q_{t-1}) = \pi_{q_1} a_{q_1 q_2} \cdots a_{q_{T-1} q_T}$$

$\pi_{q_1}$  is the probability that the first state is  $q_1$ ,  $a_{q_1 q_2}$  is the probability of going from  $q_1$  to  $q_2$ , and so on. We multiply these probabilities to get the probability of the whole sequence.

Let us now see an example (Rabiner and Juang 1986) to help us demonstrate. Assume we have  $N$  urns where each urn contains balls of only one color. So there is an urn of red balls, another of blue balls, and so forth.

Somebody draws balls from urns one by one and shows us their color. Let  $q_t$  denote the color of the ball drawn at time  $t$ . Let us say we have three states:

$S_1$  : red,  $S_2$  = blue,  $S_3$  : green

with initial probabilities:

$$\mathbf{\Pi} = [0.5, 0.2, 0.3]^T$$

$a_{ij}$  is the probability of drawing from urn  $j$  (a ball of color  $j$ ) after drawing a ball of color  $i$  from urn  $i$ . The transition matrix is, for example,

$$\mathbf{A} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

Given  $\mathbf{\Pi}$  and  $\mathbf{A}$ , it is easy to generate  $K$  random sequences each of length  $T$ . Let us see how we can calculate the probability of a sequence. Assume that the first four balls are “red, red, green, green.” This corresponds to the observation sequence  $O = \{S_1, S_1, S_3, S_3\}$ . Its probability is

$$\begin{aligned} P(O|\mathbf{A}, \mathbf{\Pi}) &= P(S_1) \cdot P(S_1|S_1) \cdot P(S_3|S_1) \cdot P(S_3|S_3) \\ &= \pi_1 \cdot a_{11} \cdot a_{13} \cdot a_{33} \\ (15.7) \quad &= 0.5 \cdot 0.4 \cdot 0.3 \cdot 0.8 = 0.048 \end{aligned}$$

Now, let us see how we can learn the parameters,  $\mathbf{\Pi}, \mathbf{A}$ . Given  $K$  sequences of length  $T$ , where  $q_t^k$  is the state at time  $t$  of sequence  $k$ , the initial probability estimate is the number of sequences starting with  $S_i$  divided by the number of sequences:

$$(15.8) \quad \hat{\pi}_i = \frac{\#\{\text{sequences starting with } S_i\}}{\#\{\text{sequences}\}} = \frac{\sum_k 1(q_1^k = S_i)}{K}$$

where  $1(b)$  is 1 if  $b$  is true and 0 otherwise.

As for the transition probabilities, the estimate for  $a_{ij}$  is the number of transitions from  $S_i$  to  $S_j$  divided by the total number of transitions from  $S_i$  over all sequences:

$$(15.9) \quad \hat{a}_{ij} = \frac{\#\{\text{transitions from } S_i \text{ to } S_j\}}{\#\{\text{transitions from } S_i\}} = \frac{\sum_k \sum_{t=1}^{T-1} 1(q_t^k = S_i \text{ and } q_{t+1}^k = S_j)}{\sum_k \sum_{t=1}^{T-1} 1(q_t^k = S_i)}$$

$\hat{a}_{12}$  is the number of times a blue ball follows a red ball divided by the total number of red ball draws over all sequences.

### 15.3 Hidden Markov Models

HIDDEN MARKOV  
MODEL

In a *hidden Markov model* (HMM), the states are not observable, but when we visit a state, an observation is recorded that is a probabilistic function of the state. We assume a discrete observation in each state from the set  $\{v_1, v_2, \dots, v_M\}$ :

$$(15.10) \quad b_j(m) \equiv P(O_t = v_m | q_t = S_j)$$

OBSERVATION  
PROBABILITY  
EMISSION  
PROBABILITY

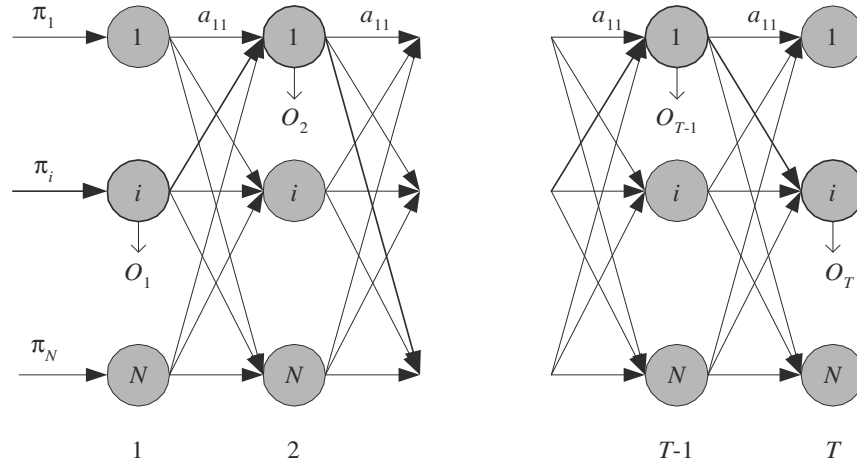
$b_j(m)$  is the *observation*, or *emission probability*, that we observe  $v_m$ ,  $m = 1, \dots, M$  in state  $S_j$ . We again assume a homogeneous model in which the probabilities do not depend on  $t$ . The values thus observed constitute the observation sequence  $O$ . The state sequence  $Q$  is not observed, that is what makes the model “hidden,” but it should be inferred from the observation sequence  $O$ . Note that there are typically many different state sequences  $Q$  that could have generated the same observation sequence  $O$ , but with different probabilities; just as, given an iid sample from a normal distribution, there are an infinite number of  $(\mu, \sigma)$  value pairs possible, we are interested in the one having the highest likelihood of generating the sample.

Note also that in this case of a hidden Markov model, there are two sources of randomness. In addition to randomly moving from one state to another, the observation in a state is also random.

Let us go back to our example. The hidden case corresponds to the urn-and-ball example where each urn contains balls of different colors. Let  $b_j(m)$  denote the probability of drawing a ball of color  $m$  from urn  $j$ . We again observe a sequence of ball colors but without knowing the sequence of urns from which the balls were drawn. So it is as if now the urns are placed behind a curtain and somebody picks a ball at random from one of the urns and shows us only the ball, without showing us the urn from which it is picked. The ball is returned to the urn to keep the probabilities the same. The number of ball colors may be different from the number of urns. For example, let us say we have three urns and the observation sequence is

$$O = \{\text{red, red, green, blue, yellow}\}$$

In the previous case, knowing the observation (ball color), we knew the state (urn) exactly because there were separate urns for separate colors and each urn contained balls of only one color. The observable model is a special case of the hidden model where  $M = N$  and  $b_j(m)$  is 1 if  $j = m$



**Figure 15.2** An HMM unfolded in time as a lattice (or trellis) showing all the possible trajectories. One path, shown in thicker lines, is the actual (unknown) state trajectory that generated the observation sequence.

and 0 otherwise. But in the case of a hidden model, a ball could have been picked from any urn. In this case, for the same observation sequence  $O$ , there may be many possible state sequences  $Q$  that could have generated  $O$  (see figure 15.2).

To summarize and formalize, an HMM has the following elements:

1.  $N$ : Number of states in the model

$$S = \{S_1, S_2, \dots, S_N\}$$

2.  $M$ : Number of distinct observation symbols in the *alphabet*

$$V = \{v_1, v_2, \dots, v_M\}$$

3. State transition probabilities:

$$\mathbf{A} = [a_{ij}] \text{ where } a_{ij} \equiv P(q_{t+1} = S_j | q_t = S_i)$$

4. Observation probabilities:

$$\mathbf{B} = [b_j(m)] \text{ where } b_j(m) \equiv P(O_t = v_m | q_t = S_j)$$

5. Initial state probabilities:

$$\mathbf{\Pi} = [\pi_i] \text{ where } \pi_i \equiv P(q_1 = S_i)$$

$N$  and  $M$  are implicitly defined in the other parameters so  $\lambda = (\mathbf{A}, \mathbf{B}, \mathbf{\Pi})$  is taken as the parameter set of an HMM. Given  $\lambda$ , the model can be used to generate an arbitrary number of observation sequences of arbitrary length, but as usual, we are interested in the other direction, that of estimating the parameters of the model given a training set of sequences.

## 15.4 Three Basic Problems of HMMs

Given a number of sequences of observations, we are interested in three problems:

1. Given a model  $\lambda$ , we would like to evaluate the probability of any given observation sequence,  $O = \{O_1 O_2 \cdots O_T\}$ , namely,  $P(O|\lambda)$ .
2. Given a model  $\lambda$  and an observation sequence  $O$ , we would like to find out the state sequence  $Q = \{q_1 q_2 \cdots q_T\}$ , which has the highest probability of generating  $O$ ; namely, we want to find  $Q^*$  that maximizes  $P(Q|O, \lambda)$ .
3. Given a training set of observation sequences,  $\mathcal{X} = \{O^k\}_k$ , we would like to learn the model that maximizes the probability of generating  $\mathcal{X}$ ; namely, we want to find  $\lambda^*$  that maximizes  $P(\mathcal{X}|\lambda)$ .

Let us see solutions to these one by one, with each solution used to solve the next problem, until we get to calculating  $\lambda$  or learning a model from data.

## 15.5 Evaluation Problem

Given an observation sequence  $O = \{O_1 O_2 \cdots O_T\}$  and a state sequence  $Q = \{q_1 q_2 \cdots q_T\}$ , the probability of observing  $O$  given the state sequence  $Q$  is simply

$$(15.11) \quad P(O|Q, \lambda) = \prod_{t=1}^T P(O_t|q_t, \lambda) = b_{q_1}(O_1) \cdot b_{q_2}(O_2) \cdots b_{q_T}(O_T)$$

which we cannot calculate because we do not know the state sequence. The probability of the state sequence  $Q$  is

$$(15.12) \quad P(Q|\lambda) = P(q_1) \prod_{t=2}^T P(q_t|q_{t-1}) = \pi_{q_1} a_{q_1 q_2} \cdots a_{q_{T-1} q_T}$$

Then the joint probability is

$$(15.13) \quad \begin{aligned} P(O, Q|\lambda) &= P(q_1) \prod_{t=2}^T P(q_t|q_{t-1}) \prod_{t=1}^T P(O_t|q_t) \\ &= \pi_{q_1} b_{q_1}(O_1) a_{q_1 q_2} b_{q_2}(O_2) \cdots a_{q_{T-1} q_T} b_{q_T}(O_T) \end{aligned}$$

We can compute  $P(O|\lambda)$  by marginalizing over the joint, namely, by summing up over all possible  $Q$ :

$$P(O|\lambda) = \sum_{\text{all possible } Q} P(O, Q|\lambda)$$

However, this is not practical since there are  $N^T$  possible  $Q$ , assuming that all the probabilities are nonzero. Fortunately, there is an efficient procedure to calculate  $P(O|\lambda)$ , which is called the *forward-backward procedure* (see figure 15.3). It is based on the idea of dividing the observation sequence into two parts: the first one starting from time 1 until time  $t$ , and the second one from time  $t + 1$  until  $T$ .

FORWARD-BACKWARD  
PROCEDURE

FORWARD VARIABLE

We define the *forward variable*  $\alpha_t(i)$  as the probability of observing the partial sequence  $\{O_1 \cdots O_t\}$  until time  $t$  and being in  $S_i$  at time  $t$ , given the model  $\lambda$ :

$$(15.14) \quad \alpha_t(i) \equiv P(O_1 \cdots O_t, q_t = S_i|\lambda)$$

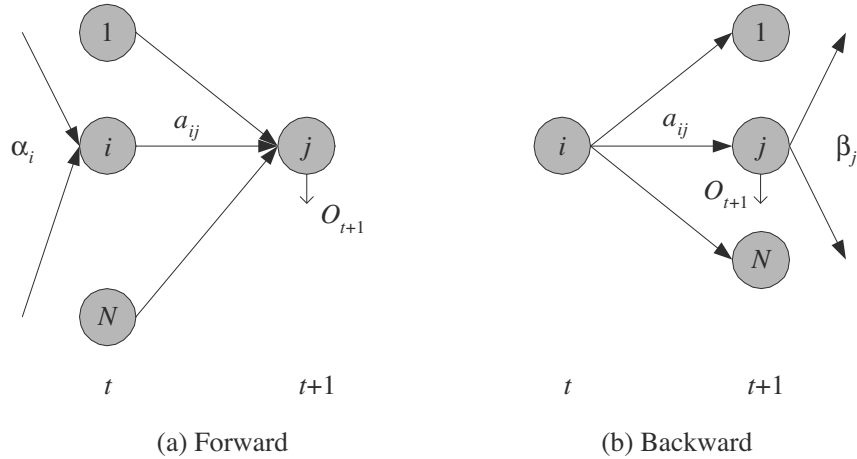
The nice thing about this is that it can be calculated recursively by accumulating results on the way.

■ Initialization:

$$(15.15) \quad \begin{aligned} \alpha_1(i) &\equiv P(O_1, q_1 = S_i|\lambda) \\ &= P(O_1|q_1 = S_i, \lambda)P(q_1 = S_i|\lambda) \\ &= \pi_i b_i(O_1) \end{aligned}$$

■ Recursion (see figure 15.3a):

$$\alpha_{t+1}(j) \equiv P(O_1 \cdots O_{t+1}, q_{t+1} = S_j|\lambda)$$



**Figure 15.3** Forward-backward procedure: (a) computation of  $\alpha_t(j)$  and (b) computation of  $\beta_t(i)$ .

$$\begin{aligned}
 &= P(O_1 \cdots O_{t+1} | q_{t+1} = S_j, \lambda) P(q_{t+1} = S_j | \lambda) \\
 &= P(O_1 \cdots O_t | q_{t+1} = S_j, \lambda) P(O_{t+1} | q_{t+1} = S_j, \lambda) P(q_{t+1} = S_j | \lambda) \\
 &= P(O_1 \cdots O_t, q_{t+1} = S_j | \lambda) P(O_{t+1} | q_{t+1} = S_j, \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \sum_i P(O_1 \cdots O_t, q_t = S_i, q_{t+1} = S_j | \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \\
 &\quad \sum_i P(O_1 \cdots O_t, q_{t+1} = S_j | q_t = S_i, \lambda) P(q_t = S_i | \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \\
 &\quad \sum_i P(O_1 \cdots O_t | q_t = S_i, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) P(q_t = S_i | \lambda) \\
 &= P(O_{t+1} | q_{t+1} = S_j, \lambda) \\
 &\quad \sum_i P(O_1 \cdots O_t, q_t = S_i | \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\
 (15.16) \quad &= \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1})
 \end{aligned}$$

$\alpha_t(i)$  explains the first  $t$  observations and ends in state  $S_i$ . We multiply this by the probability  $a_{ij}$  to move to state  $S_j$ , and because there are

$N$  possible previous states, we need to sum up over all such possible previous  $S_i$ .  $b_j(O_{t+1})$  then is the probability we generate the  $(t + 1)$ st observation while in state  $S_j$  at time  $t + 1$ .

When we calculate the forward variables, it is easy to calculate the probability of the observation sequence:

$$\begin{aligned}
 P(O|\lambda) &= \sum_{i=1}^N P(O, q_T = S_i|\lambda) \\
 (15.17) \quad &= \sum_{i=1}^N \alpha_T(i)
 \end{aligned}$$

$\alpha_T(i)$  is the probability of generating the full observation sequence and ending up in state  $S_i$ . We need to sum up over all such possible final states.

Computing  $\alpha_t(i)$  is  $\mathcal{O}(N^2T)$ , and this solves our first evaluation problem in a reasonable amount of time. We do not need it now but let us similarly define the *backward variable*,  $\beta_t(i)$ , which is the probability of being in  $S_i$  at time  $t$  and observing the partial sequence  $O_{t+1} \cdots O_T$ :

BACKWARD VARIABLE

$$(15.18) \quad \beta_t(i) \equiv P(O_{t+1} \cdots O_T | q_t = S_i, \lambda)$$

This can again be recursively computed as follows, this time going in the backward direction:

- Initialization (arbitrarily to 1):

$$\beta_T(i) = 1$$

- Recursion (see figure 15.3b):

$$\begin{aligned}
 \beta_t(i) &\equiv P(O_{t+1} \cdots O_T | q_t = S_i, \lambda) \\
 &= \sum_j P(O_{t+1} \cdots O_T, q_{t+1} = S_j | q_t = S_i, \lambda) \\
 &= \sum_j P(O_{t+1} \cdots O_T | q_{t+1} = S_j, q_t = S_i, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\
 &= \sum_j P(O_{t+1} | q_{t+1} = S_j, q_t = S_i, \lambda) \\
 &\quad P(O_{t+2} \cdots O_T | q_{t+1} = S_j, q_t = S_i, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\
 &= \sum_j P(O_{t+1} | q_{t+1} = S_j, \lambda)
 \end{aligned}$$

$$\begin{aligned}
 & P(O_{t+2} \cdots O_T | q_{t+1} = S_j, \lambda) P(q_{t+1} = S_j | q_t = S_i, \lambda) \\
 (15.19) \quad & = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)
 \end{aligned}$$

When in state  $S_i$ , we can go to  $N$  possible next states  $S_j$ , each with probability  $a_{ij}$ . While there, we generate the  $(t + 1)$ st observation and  $\beta_{t+1}(j)$  explains all the observations after time  $t + 1$ , continuing from there.

One word of caution about implementation is necessary here: Both  $\alpha_t$  and  $\beta_t$  values are calculated by multiplying small probabilities, and with long sequences we risk getting underflow. To avoid this, at each time step, we normalize  $\alpha_t(i)$  by multiplying it with

$$c_t = \frac{1}{\sum_j \alpha_t(j)}$$

We also normalize  $\beta_t(i)$  by multiplying it with the same  $c_t$  ( $\beta_t(i)$  do not sum to 1). We cannot use equation 15.17 after normalization; instead, we have (Rabiner 1989)

$$(15.20) \quad P(O|\lambda) = \frac{1}{\prod_t c_t} \text{ or } \log P(O|\lambda) = - \sum_t \log c_t$$

## 15.6 Finding the State Sequence

We now move on to the second problem, that of finding the state sequence  $Q = \{q_1 q_2 \cdots q_T\}$  having the highest probability of generating the observation sequence  $O = \{O_1 O_2 \cdots O_T\}$ , given the model  $\lambda$ .

Let us define  $y_t(i)$  as the probability of being in state  $S_i$  at time  $t$ , given  $O$  and  $\lambda$ , which can be computed as

$$\begin{aligned}
 (15.21) \quad y_t(i) & \equiv P(q_t = S_i | O, \lambda) \\
 & = \frac{P(O | q_t = S_i, \lambda) P(q_t = S_i | \lambda)}{P(O | \lambda)} \\
 & = \frac{P(O_1 \cdots O_t | q_t = S_i, \lambda) P(O_{t+1} \cdots O_T | q_t = S_i, \lambda) P(q_t = S_i | \lambda)}{\sum_{j=1}^N P(O, q_t = S_j | \lambda)} \\
 & = \frac{P(O_1 \cdots O_t, q_t = S_i | \lambda) P(O_{t+1} \cdots O_T | q_t = S_i, \lambda)}{\sum_{j=1}^N P(O | q_t = S_j, \lambda) P(q_t = S_j | \lambda)} \\
 (15.22) \quad & = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)}
 \end{aligned}$$

Here we see how nicely  $\alpha_t(i)$  and  $\beta_t(i)$  split the sequence between them: the forward variable  $\alpha_t(i)$  explains the starting part of the sequence until time  $t$  and ends in  $S_i$ , and the backward variable  $\beta_t(i)$  takes it from there and explains the ending part until time  $T$ .

The numerator  $\alpha_t(i)\beta_t(i)$  explains the whole sequence given that at time  $t$ , the system is in state  $S_i$ . We need to normalize by dividing this over all possible intermediate states that can be traversed at time  $t$ , and guarantee that  $\sum_i \gamma_t(i) = 1$ .

To find the state sequence, for each time step  $t$ , we can choose the state that has the highest probability:

$$(15.23) \quad q_t^* = \arg \max_i \gamma_t(i)$$

but this may choose  $S_i$  and  $S_j$  as the most probable states at time  $t$  and  $t + 1$  even when  $a_{ij} = 0$ . To find the single best state *sequence* (path), we use the *Viterbi algorithm*, based on dynamic programming, which takes such transition probabilities into account.

VITERBI ALGORITHM

Given state sequence  $Q = q_1 q_2 \cdots q_T$  and observation sequence  $O = O_1 \cdots O_T$ , we define  $\delta_t(i)$  as the probability of the highest probability path at time  $t$  that accounts for the first  $t$  observations and ends in  $S_i$ :

$$(15.24) \quad \delta_t(i) \equiv \max_{q_1 q_2 \cdots q_{t-1}} p(q_1 q_2 \cdots q_{t-1}, q_t = S_i, O_1 \cdots O_t | \lambda)$$

Then we can recursively calculate  $\delta_{t+1}(i)$  and the optimal path can be read by backtracking from  $T$ , choosing the most probable at each instant. The algorithm is as follows:

1. Initialization:

$$\begin{aligned} \delta_1(i) &= \pi_i b_i(O_1) \\ \psi_1(i) &= 0 \end{aligned}$$

2. Recursion:

$$\begin{aligned} \delta_t(j) &= \max_i \delta_{t-1}(i) a_{ij} \cdot b_j(O_t) \\ \psi_t(j) &= \arg \max_i \delta_{t-1}(i) a_{ij} \end{aligned}$$

3. Termination:

$$\begin{aligned} p^* &= \max_i \delta_T(i) \\ q_T^* &= \arg \max_i \delta_T(i) \end{aligned}$$

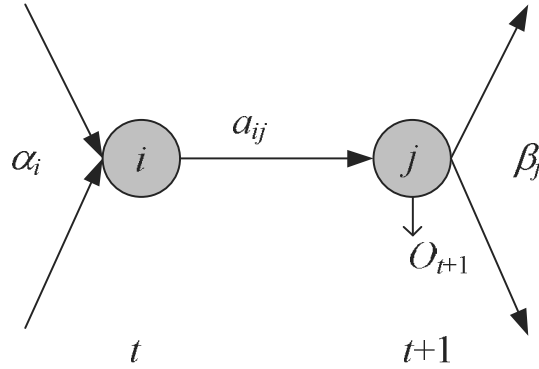


Figure 15.4 Computation of arc probabilities,  $\xi_t(i, j)$ .

4. Path (state sequence) backtracking:

$$q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T - 1, T - 2, \dots, 1$$

Using the lattice structure of figure 15.2,  $\psi_t(j)$  keeps track of the state that maximizes  $\delta_t(j)$  at time  $t - 1$ , that is, the best previous state. The Viterbi algorithm has the same complexity with the forward phase, where instead of the sum, we take the maximum at each step.

### 15.7 Learning Model Parameters

We now move on to the third problem, learning an HMM from data. The approach is maximum likelihood, and we would like to calculate  $\lambda^*$  that maximizes the likelihood of the sample of training sequences,  $\mathcal{X} = \{O^k\}_{k=1}^K$ , namely,  $P(\mathcal{X}|\lambda)$ . We start by defining a new variable that will become handy later on.

We define  $\xi_t(i, j)$  as the probability of being in  $S_i$  at time  $t$  and in  $S_j$  at time  $t + 1$ , given the whole observation  $O$  and  $\lambda$ :

$$(15.25) \quad \xi_t(i, j) \equiv P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

which can be computed as (see figure 15.4)

$$\begin{aligned} \xi_t(i, j) &\equiv P(q_t = S_i, q_{t+1} = S_j | O, \lambda) \\ &= \frac{P(O | q_t = S_i, q_{t+1} = S_j, \lambda) P(q_t = S_i, q_{t+1} = S_j | \lambda)}{P(O | \lambda)} \end{aligned}$$

$$\begin{aligned}
&= \frac{P(O|q_t = S_i, q_{t+1} = S_j, \lambda)P(q_{t+1} = S_j|q_t = S_i, \lambda)P(q_t = S_i|\lambda)}{P(O|\lambda)} \\
&= \left(\frac{1}{P(O|\lambda)}\right)P(O_1 \cdots O_t|q_t = S_i, \lambda)P(O_{t+1}|q_{t+1} = S_j, \lambda) \\
&\quad P(O_{t+2} \cdots O_T|q_{t+1} = S_j, \lambda)a_{ij}P(q_t = S_i|\lambda) \\
&= \left(\frac{1}{P(O|\lambda)}\right)P(O_1 \cdots O_t, q_t = S_i|\lambda)P(O_{t+1}|q_{t+1} = S_j, \lambda) \\
&\quad P(O_{t+2} \cdots O_T|q_{t+1} = S_j, \lambda)a_{ij} \\
&= \frac{\alpha_t(i)b_j(O_{t+1})\beta_{t+1}(j)a_{ij}}{\sum_k \sum_l P(q_t = S_k, q_{t+1} = S_l, O|\lambda)} \\
(15.26) \quad &= \frac{\alpha_t(i)a_{ij}b_j(O_{t+1})\beta_{t+1}(j)}{\sum_k \sum_l \alpha_t(k)a_{kl}b_l(O_{t+1})\beta_{t+1}(l)}
\end{aligned}$$

$\alpha_t(i)$  explains the first  $t$  observations and ends in state  $S_i$  at time  $t$ . We move on to state  $S_j$  with probability  $a_{ij}$ , generate the  $(t+1)$ st observation, and continue from  $S_j$  at time  $t+1$  to generate the rest of the observation sequence. We normalize by dividing for all such possible pairs that can be visited at time  $t$  and  $t+1$ .

If we want, we can also calculate the probability of being in state  $S_i$  at time  $t$  by marginalizing over the arc probabilities for all possible next states:

$$(15.27) \quad \gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

SOFT COUNTS

Note that if the Markov model were not hidden but observable, both  $\gamma_t(i)$  and  $\xi_t(i, j)$  would be 0/1. In this case when they are not, we estimate them with posterior probabilities that give us *soft counts*. This is just like the difference between supervised classification and unsupervised clustering where we did and did not know the class labels, respectively. In unsupervised clustering using EM (section 7.4), not knowing the class labels, we estimated them first (in the E-step) and calculated the parameters with these estimates (in the M-step).

BAUM-WELCH  
ALGORITHM

Similarly here we have the *Baum-Welch algorithm*, which is an EM procedure. At each iteration, first in the E-step, we compute  $\xi_t(i, j)$  and  $\gamma_t(i)$  values given the current  $\lambda = (\mathbf{A}, \mathbf{B}, \mathbf{\Pi})$ , and then in the M-step, we recalculate  $\lambda$  given  $\xi_t(i, j)$  and  $\gamma_t(i)$ . These two steps are alternated until convergence during which, it has been shown,  $P(O|\lambda)$  never decreases.

Assume indicator variables  $z_i^t$  as

$$(15.28) \quad z_i^t = \begin{cases} 1 & \text{if } q_t = S_i \\ 0 & \text{otherwise} \end{cases}$$

and

$$(15.29) \quad z_{ij}^t = \begin{cases} 1 & \text{if } q_t = S_i \text{ and } q_{t+1} = S_j \\ 0 & \text{otherwise} \end{cases}$$

These are 0/1 in the case of an observable Markov model and are hidden random variables in the case of an HMM. In this latter case, we estimate them in the E-step as

$$(15.30) \quad \begin{aligned} E[z_i^t] &= \gamma_t(i) \\ E[z_{ij}^t] &= \xi_t(i, j) \end{aligned}$$

In the M-step, we calculate the parameters given these estimated values. The expected number of transitions from  $S_i$  to  $S_j$  is  $\sum_t \xi_t(i, j)$  and the total number of transitions from  $S_i$  is  $\sum_t \gamma_t(i)$ . The ratio of these two gives us the probability of transition from  $S_i$  to  $S_j$  at any time:

$$(15.31) \quad \hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

Note that this is the same as equation 15.9, except that the actual counts are replaced by estimated soft counts.

The probability of observing  $v_m$  in  $S_j$  is the expected number of times  $v_m$  is observed when the system is in  $S_j$  over the total number of times the system is in  $S_j$ :

$$(15.32) \quad \hat{b}_j(m) = \frac{\sum_{t=1}^T \gamma_t(j) 1(O_t = v_m)}{\sum_{t=1}^T \gamma_t(j)}$$

When there are multiple observation sequences

$$\mathcal{X} = \{O^k\}_{k=1}^K$$

which we assume to be independent

$$P(\mathcal{X}|\lambda) = \prod_{k=1}^K P(O^k|\lambda)$$

the parameters are now averages over all observations in all sequences:

$$(15.33) \quad \begin{aligned} \hat{a}_{ij} &= \frac{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \xi_t^k(i, j)}{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \gamma_t^k(i)} \\ \hat{b}_j(m) &= \frac{\sum_{k=1}^K \sum_{t=1}^{T_k} \gamma_t^k(j) \mathbf{1}(O_t^k = v_m)}{\sum_{k=1}^K \sum_{t=1}^{T_k} \gamma_t^k(j)} \\ \hat{\pi}_i &= \frac{\sum_{k=1}^K \gamma_1^k(i)}{K} \end{aligned}$$

## 15.8 Continuous Observations

In our discussion, we assumed discrete observations modeled as a multinomial

$$(15.34) \quad P(O_t | q_t = S_j, \lambda) = \prod_{m=1}^M b_j(m)^{r_m^t}$$

where

$$(15.35) \quad r_m^t = \begin{cases} 1 & \text{if } O_t = v_m \\ 0 & \text{otherwise} \end{cases}$$

If the inputs are continuous, one possibility is to discretize them and then use these discrete values as observations. Typically, a vector quantizer (section 7.3) is used for this purpose of converting continuous values to the discrete index of the closest reference vector. For example, in speech recognition, a word utterance is divided into short speech segments corresponding to phonemes or part of phonemes; after preprocessing, these are discretized using a vector quantizer and an HMM is then used to model a word utterance as a sequence of them.

We remember that  $k$ -means used for vector quantization is the hard version of a Gaussian mixture model:

$$(15.36) \quad p(O_t | q_t = S_j, \lambda) = \sum_{l=1}^L P(G_l) p(O_t | q_t = S_j, G_l, \lambda)$$

where

$$(15.37) \quad p(O_t | q_t = S_j, G_l, \lambda) \sim \mathcal{N}(\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)$$

and the observations are kept continuous. In this case of Gaussian mixtures, EM equations can be derived for the component parameters (with

suitable regularization to keep the number of parameters in check) and the mixture proportions (Rabiner 1989).

Let us see the case of a scalar continuous observation,  $O_t \in \mathfrak{R}$ . The easiest is to assume a normal distribution:

$$(15.38) \quad p(O_t | q_t = S_j, \lambda) \sim \mathcal{N}(\mu_j, \sigma_j^2)$$

which implies that in state  $S_j$ , the observation is drawn from a normal with mean  $\mu_j$  and variance  $\sigma_j^2$ . The M-step equations in this case are

$$(15.39) \quad \begin{aligned} \hat{\mu}_j &= \frac{\sum_t \gamma_t(j) O_t}{\sum_t \gamma_t(j)} \\ \hat{\sigma}_j^2 &= \frac{\sum_t \gamma_t(j) (O_t - \hat{\mu}_j)^2}{\sum_t \gamma_t(j)} \end{aligned}$$

## 15.9 The HMM with Input

In some applications, additional to the observation sequence  $O_t$ , we have an input sequence,  $x_t$ . We can condition the observation  $O_t$  in state  $S_j$  on the input  $x^t$ , and write  $P(O_t | q_t = S_j, x_t)$ . In the case when the observations are continuous scalars, we replace equation 15.38 with a generalized model

$$(15.40) \quad p(O_t | q_t = S_j, x_t, \lambda) \sim \mathcal{N}(g_j(x^t | \theta_j), \sigma_j^2)$$

where, for example, assuming a linear model, we have

$$(15.41) \quad g_j(x^t | w_j, w_{j0}) = w_j x^t + w_{j0}$$

If the observations are discrete and multinomial, we have a classifier taking  $x^t$  as input and generating a 1-of- $M$  output, or we can generate posterior class probabilities and keep the observations continuous.

Similarly, the state transition probabilities can also be conditioned on the input, namely,  $P(q_{t+1} = S_j | q_t = S_i, x_t)$ , which is implemented by a classifier choosing the state at time  $t + 1$  as a function of the state at time  $t$  and the input. This is a *Markov mixture of experts* (Meila and Jordan 1996) and is a generalization of the mixture of experts architecture (section 12.8) where the gating network keeps track of the decision it made in the previous time step. Such an architecture is also called an *input-output HMM* (Bengio and Frasconi 1996) and has the advantage that the model is no longer homogeneous; different observation and transition

MARKOV MIXTURE OF  
EXPERTS

INPUT-OUTPUT HMM

probabilities are used at different time steps. There is still a single model for each state, parameterized by  $\theta_j$ , but it generates different transition or observation probabilities depending on the input seen. It is possible that the input is not a single value but a window around time  $t$  making the input a vector; this allows handling applications where the input and observation sequences have different lengths.

Even if there is no other explicit input sequence, an HMM with input can be used by generating an “input” through some prespecified function of previous observations

$$\mathbf{x}_t = \mathbf{f}(O_{t-\tau}, \dots, O_{t-1})$$

thereby providing a window of size  $\tau$  of contextual input.

### 15.10 Model Selection in HMM

Just like any model, the complexity of an HMM should be tuned so as to balance its complexity with the size and properties of the data at hand. One possibility is to tune the topology of the HMM. In a fully connected (ergodic) HMM, there is transition from a state to any other state, which makes  $\mathbf{A}$  a full  $N \times N$  matrix. In some applications, only certain transitions are allowed, with the disallowed transitions having their  $a_{ij} = 0$ . When there are fewer possible next states,  $N' < N$ , the complexity of forward-backward passes and the Viterbi procedure is  $\mathcal{O}(NN'T)$  instead of  $\mathcal{O}(N^2T)$ .

LEFT-TO-RIGHT HMMs

For example, in speech recognition, *left-to-right HMMs* are used, which have their states ordered in time so that as time increases, the state index increases or stays the same. Such a constraint allows modeling sequences whose properties change over time as in speech, and when we get to a state, we know approximately the states preceding it. There is the property that we never move to a state with a smaller index, namely,  $a_{ij} = 0$ , for  $j < i$ . Large changes in state indices are not allowed either, namely,  $a_{ij} = 0$ , for  $j > i + \tau$ . The example of the left-to-right HMM given in figure 15.5 with  $\tau = 2$  has the state transition matrix

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & 0 \\ 0 & a_{22} & a_{23} & a_{24} \\ 0 & 0 & a_{33} & a_{34} \\ 0 & 0 & 0 & a_{44} \end{bmatrix}$$

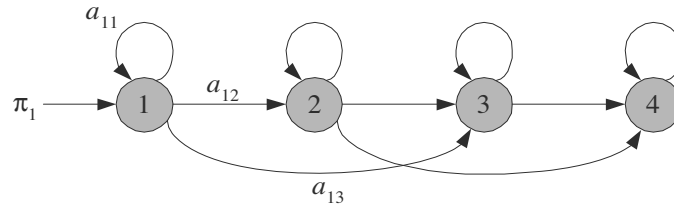


Figure 15.5 Example of a left-to-right HMM.

Another factor that determines the complexity of an HMM is the number of states  $N$ . Because the states are hidden, their number is not known and should be chosen before training. This is determined using prior information and can be fine-tuned by cross-validation, namely, by checking the likelihood of validation sequences.

When used for classification, we have a set of HMMs, each one modeling the sequences belonging to one class. For example, in spoken word recognition, examples of each word train a separate model,  $\lambda_i$ . Given a new word utterance  $O$  to classify, all of the separate word models are evaluated to calculate  $P(O|\lambda_i)$ . We then use Bayes' rule to get the posterior probabilities

$$(15.42) \quad P(\lambda_i|O) = \frac{P(O|\lambda_i)P(\lambda_i)}{\sum_j P(O|\lambda_j)P(\lambda_j)}$$

where  $P(\lambda_i)$  is the prior probability of word  $i$ . The utterance is assigned to the word having the highest posterior. This is the likelihood-based approach; there is also work on discriminative HMM trained directly to maximize the posterior probabilities. When there are several pronunciations of the same word, these are defined as parallel paths in the HMM for the word.

PHONES In the case of a continuous input like speech, the difficult task is that of segmenting the signal into small discrete observations. Typically, *phones* are used that are taken as the primitive parts, and combining them, longer sequences (e.g., words) are formed. Each phone is recognized in parallel (by the vector quantizer), then the HMM is used to combine them serially. If the speech primitives are simple, then the HMM becomes complex and vice versa. In connected speech recognition where the words are not uttered one by one with clear pauses between them, there is a hierarchy of HMMs at several levels; one combines phones to recognize words,

another combines words to recognize sentences by building a language model, and so forth.

Hybrid neural network/HMM models were also used for speech recognition (Morgan and Bourlard 1995). In such a model, a multilayer perceptron (chapter 11) is used to capture temporally local but possibly complex and nonlinear primitives, for example, phones, while the HMM is used to learn the temporal structure. The neural network acts as a preprocessor and translates the raw observations in a time window to a form that is easier to model than the output of a vector quantizer.

An HMM can be visualized as a graphical model and evaluation in an HMM is a special case of the belief propagation algorithm, as we will see in chapter 16. The reason that we devote a special chapter is the widespread successful use of this particular model, especially in automatic speech recognition. When we discuss graphical models in detail, we will see how the basic HMM architecture can be extended—for example, by having multiple sequences, or by introducing hidden (latent) variables that can simplify the model.

### 15.11 Notes

The HMM is a mature technology, and there are HMM-based commercial speech recognition systems in actual use (Rabiner and Juang 1993; Jelinek 1997). In section 11.12, we discussed how to train multilayer perceptrons for recognizing sequences. HMMs have the advantage over time delay neural networks in that no time window needs to be defined a priori, and they train better than recurrent neural networks. HMMs are applied to diverse sequence recognition tasks. Applications of HMMs to bioinformatics is given in Baldi and Brunak 1998, and to natural language processing in Manning and Schütze 1999. It is also applied to online handwritten character recognition, which differs from optical recognition in that the writer writes on a touch-sensitive pad and the input is a sequence of  $(x, y)$  coordinates of the pen tip as it moves over the pad and is not a static image. Bengio et al. (1995) explain a hybrid system for online recognition where an MLP recognizes individual characters, and an HMM combines them to recognize words. Various applications of the HMM and several extensions, for example, discriminative HMMs, are discussed in Bengio 1999. A more recent survey of what HMMs can and cannot do is Bilmes 2006.

In any such recognition system, one critical point is to decide how much to do things in parallel and what to leave to serial processing. In speech recognition, phonemes may be recognized by a parallel system that corresponds to assuming that all the phoneme sound is uttered in one time step. The word is then recognized serially by combining the phonemes. In an alternative system, phonemes themselves may be designed as a sequence of simpler speech sounds, if the same phoneme has many versions, for example, depending on the previous and following phonemes. Doing things in parallel is good but only to a degree; one should find the ideal balance of parallel and serial processing. To be able to call anyone at the touch of a button, we would need millions of buttons on our telephone; instead, we have ten buttons and we press them in a sequence to dial the number.

We will discuss graphical models in chapter 16 where we will see that HMMs can be considered a special class of graphical models and inference and learning operations on HMMs are analogous to their counterparts in Bayesian networks (Smyth, Heckerman, and Jordan 1997). As we will see shortly, there are various extensions to HMMs like *factorial HMMs* where at each time step, there are a number of states that collectively generate the observation and *tree-structured HMMs* where there is a hierarchy of states. The general formalism also allows us to treat continuous as well as discrete states, known as *linear dynamical systems*. For some of these models, exact inference is not possible and one needs to use approximation or sampling methods (Ghahramani 2001).

## 15.12 Exercises

1. Given the observable Markov model with three states,  $S_1, S_2, S_3$ , initial probabilities

$$\boldsymbol{\Pi} = [0.5, 0.2, 0.3]^T$$

and transition probabilities

$$\mathbf{A} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

generate 100 sequences of 1,000 states.

2. Using the data generated by the previous exercise, estimate  $\boldsymbol{\Pi}, \mathbf{A}$  and compare with the parameters used to generate the data.

3. Formalize a second-order Markov model. What are the parameters? How can we calculate the probability of a given state sequence? How can the parameters be learned for the case of an observable model?
4. Show that any second- (or higher-order) Markov model can be converted to a first-order Markov model.
5. Some researchers define a Markov model as generating an observation while traversing an arc, instead of on arrival at a state. Is this model any more powerful than what we have discussed?
6. Generate training and validation sequences from an HMM of your choosing. Then train different HMMs by varying the number of hidden states on the same training set and calculate the validation likelihoods. Observe how the validation likelihood changes as the number of states increases.
7. If in equation 15.38 we have multivariate observations, what will be the M-step equations?
8. Consider the urn-and-ball example where we draw *without replacement*. How will it be different?
9. Let us say at any time we have two observations from two different alphabets; for example, let us say we are observing the values of two currencies every day. How can we implement this using HMM?
10. How can we have an incremental HMM where we add new hidden states when necessary?

### 15.13 References

- Baldi, P., and S. Brunak. 1998. *Bioinformatics: The Machine Learning Approach*. Cambridge, MA: MIT Press.
- Bengio, Y. 1999. "Markovian Models for Sequential Data." *Neural Computing Surveys* 2: 129-162.
- Bengio, Y., and P. Frasconi. 1996. "Input-Output HMMs for Sequence Processing." *IEEE Transactions on Neural Networks* 7: 1231-1249.
- Bengio, Y., Y. Le Cun, C. Nohl, and C. Burges. 1995. "LeRec: A NN/HMM Hybrid for On-line Handwriting Recognition." *Neural Computation* 7: 1289-1303.
- Bilmes, J. A. 2006. "What HMMs Can Do." *IEICE Transactions on Information and Systems* E89-D: 869-891.
- Ghahramani, Z. 2001. "An Introduction to Hidden Markov Models and Bayesian Networks." *International Journal of Pattern Recognition and Artificial Intelligence* 15: 9-42.

- Jelinek, F. 1997. *Statistical Methods for Speech Recognition*. Cambridge, MA: MIT Press.
- Manning, C. D., and H. Schütze. 1999. *Foundations of Statistical Natural Language Processing*. Cambridge, MA: MIT Press.
- Meila, M., and M. I. Jordan. 1996. "Learning Fine Motion by Markov Mixtures of Experts." In *Advances in Neural Information Processing Systems 8*, ed. D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, 1003-1009. Cambridge, MA: MIT Press.
- Morgan, N., and H. Bourlard. 1995. "Continuous Speech Recognition: An Introduction to the Hybrid HMM/Connectionist Approach." *IEEE Signal Processing Magazine* 12: 25-42.
- Smyth, P., D. Heckerman, and M. I. Jordan. 1997. "Probabilistic Independence Networks for Hidden Markov Probability Models." *Neural Computation* 9: 227-269.
- Rabiner, L. R. 1989. "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition." *Proceedings of the IEEE* 77: 257-286.
- Rabiner, L. R., and B. H. Juang. 1986. "An Introduction to Hidden Markov Models." *IEEE Acoustics, Speech, and Signal Processing Magazine* 3: 4-16.
- Rabiner, L. R., and B. H. Juang. 1993. *Fundamentals of Speech Recognition*. New York: Prentice Hall.



# 16 *Graphical Models*

*Graphical models represent the interaction between variables visually and have the advantage that inference over a large number of variables can be decomposed into a set of local calculations involving a small number of variables making use of conditional independencies. After some examples of inference by hand, we discuss the concept of  $d$ -separation and the belief propagation algorithm on a variety of graphs.*

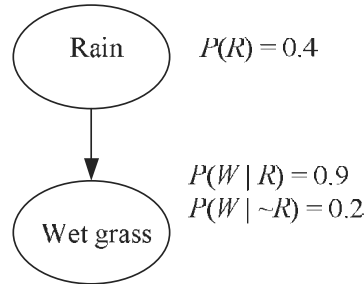
## 16.1 Introduction

GRAPHICAL MODELS  
BAYESIAN NETWORKS  
BELIEF NETWORKS  
PROBABILISTIC  
NETWORKS

DIRECTED ACYCLIC  
GRAPH

*Graphical models, also called Bayesian networks, belief networks, or probabilistic networks, are composed of nodes and arcs between the nodes. Each node corresponds to a random variable,  $X$ , and has a value corresponding to the probability of the random variable,  $P(X)$ . If there is a directed arc from node  $X$  to node  $Y$ , this indicates that  $X$  has a *direct influence* on  $Y$ . This influence is specified by the conditional probability  $P(Y|X)$ . The network is a *directed acyclic graph* (DAG); namely, there are no cycles. The nodes and the arcs between the nodes define the *structure* of the network, and the conditional probabilities are the *parameters* given the structure.*

A simple example is given in figure 16.1, which models that rain causes the grass to get wet. It rains on 40 percent of the days and when it rains, there is a 90 percent chance that the grass gets wet; maybe 10 percent of the time it does not rain long enough for us to really consider the grass wet enough. The random variables in this example are binary; they are either true or false. There is a 20 percent probability that the grass gets wet without its actually raining, for example, when a sprinkler is used.



**Figure 16.1** Bayesian network modeling that rain is the cause of wet grass.

We see that these three values completely specify the joint distribution of  $P(R, W)$ . If  $P(R) = 0.4$ , then  $P(\sim R) = 0.6$ , and similarly  $P(\sim W|R) = 0.1$  and  $P(\sim W|\sim R) = 0.8$ . The joint is written as

$$P(R, W) = P(R)P(W|R)$$

We can calculate the individual (marginal) probability of wet grass by summing up over the possible values that its parent node can take:

$$\begin{aligned} P(W) &= \sum_R P(R, W) = P(W|R)P(R) + P(W|\sim R)P(\sim R) \\ &= 0.9 \cdot 0.4 + 0.2 \cdot 0.6 = 0.48 \end{aligned}$$

If we knew that it rained, the probability of wet grass would be 0.9; if we knew for sure that it did not, it would be as low as 0.2; not knowing whether it rained or not, the probability is 0.48.

CAUSAL GRAPH

Figure 16.1 shows a *causal graph* in that it explains that the cause of wet grass is rain. Bayes' rule allows us to invert the dependencies and have a *diagnosis*. For example, knowing that the grass is wet, the probability that it rained can be calculated as follows:

$$P(R|W) = \frac{P(W|R)P(R)}{P(W)} = 0.75$$

Knowing that the grass is wet increased the probability of rain from 0.4 to 0.75; this is because  $P(W|R)$  is high and  $P(W|\sim R)$  is low.

INDEPENDENCE

We form graphs by adding nodes and arcs and generate dependencies.  $X$  and  $Y$  are *independent events* if

$$(16.1) \quad p(X, Y) = P(X)P(Y)$$

CONDITIONAL  
INDEPENDENCE

$X$  and  $Y$  are *conditionally independent events* given a third event  $Z$  if

$$(16.2) \quad P(X, Y|Z) = P(X|Z)P(Y|Z)$$

which can also be rewritten as

$$(16.3) \quad P(X|Y, Z) = P(X|Z)$$

In a graphical model, not all nodes are connected; actually, in general, a node is connected to only a small number of other nodes. Certain subgraphs imply conditional independence statements, and these allow us to break down a complex graph into smaller subsets in which inferences can be done locally and whose results are later propagated over the graph. There are three canonical cases and larger graphs are constructed using these as subgraphs.

## 16.2 Canonical Cases for Conditional Independence

### Case 1: Head-to-tail Connection

Three events may be connected serially, as seen in figure 16.2a. We see here that  $X$  and  $Z$  are independent given  $Y$ : Knowing  $Y$  tells  $Z$  everything; knowing the state of  $X$  does not add any extra knowledge for  $Z$ ; we write  $P(Z|Y, X) = P(Z|Y)$ . We say that  $Y$  *blocks* the path from  $X$  to  $Z$ , or in other words, it *separates* them in the sense that if  $Y$  is removed, there is no path between  $X$  to  $Z$ . In this case, the joint is written as

$$(16.4) \quad P(X, Y, Z) = P(X)P(Y|X)P(Z|Y)$$

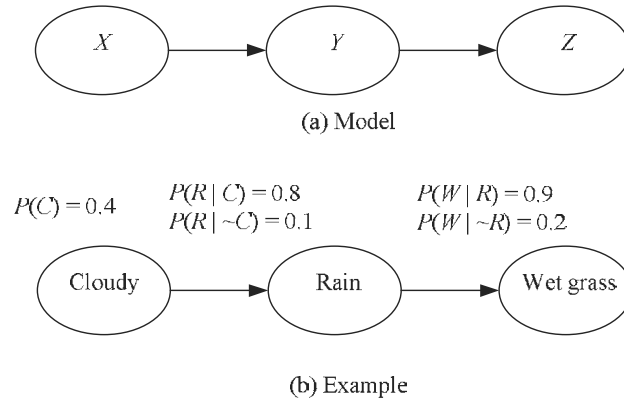
Writing the joint this way implies independence:

$$(16.5) \quad P(Z|X, Y) = \frac{P(X, Y, Z)}{P(X, Y)} = \frac{P(X)P(Y|X)P(Z|Y)}{P(X)P(Y|X)} = P(Z|Y)$$

Typically,  $X$  is the cause of  $Y$  and  $Y$  is the cause of  $Z$ . For example, as seen in figure 16.2b,  $X$  can be cloudy sky,  $Y$  can be rain, and  $Z$  can be wet grass. We can propagate information along the chain. If we do not know the state of cloudy, we have

$$\begin{aligned} P(R) &= P(R|C)P(C) + P(R|\sim C)P(\sim C) = 0.38 \\ P(W) &= P(W|R)P(R) + P(W|\sim R)P(\sim R) = 0.47 \end{aligned}$$

Let us say, in the morning we see that the weather is cloudy; what can we say about the probability that the grass will be wet? To do this, we



**Figure 16.2** Head-to-tail connection. (a) Three nodes are connected serially.  $X$  and  $Z$  are independent given the intermediate node  $Y$ :  $P(Z|Y, X) = P(Z|Y)$ . (b) Example: Cloudy weather causes rain, which in turn causes wet grass.

need to propagate evidence first to the intermediate node  $R$ , and then to the query node  $W$ .

$$P(W|C) = P(W|R)P(R|C) + P(W|\sim R)P(\sim R|C) = 0.76$$

Knowing that the weather is cloudy increased the probability of wet grass. We can also propagate evidence back using Bayes' rule. Let us say that we were traveling and on our return, see that our grass is wet; what is the probability that the weather was cloudy that day? We use Bayes' rule to invert the direction:

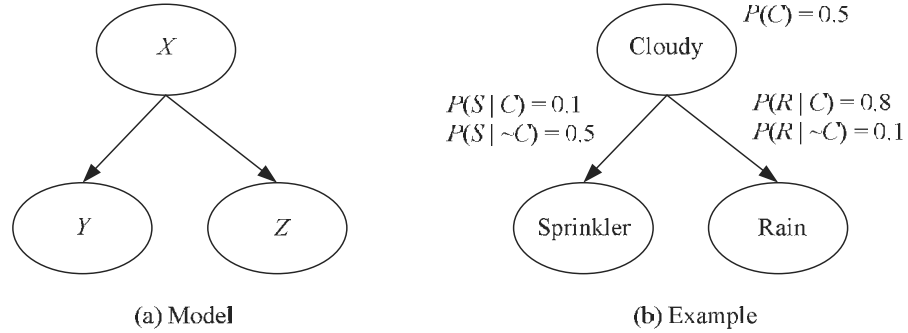
$$P(C|W) = \frac{P(W|C)P(C)}{P(W)} = 0.65$$

Knowing that the grass is wet increased the probability of cloudy weather from its default (prior) value of 0.4 to 0.65.

### Case 2: Tail-to-tail Connection

$X$  may be the parent of two nodes  $Y$  and  $Z$ , as shown in figure 16.3a. The joint density is written as

$$(16.6) \quad P(X, Y, Z) = P(X)P(Y|X)P(Z|X)$$



**Figure 16.3** Tail-to-tail connection.  $X$  is the parent of two nodes  $Y$  and  $Z$ . The two child nodes are independent given the parent:  $P(Y|X, Z) = P(Y|X)$ . In the example, cloudy weather causes rain and also makes us less likely to turn the sprinkler on.

Normally  $Y$  and  $Z$  are dependent through  $X$ ; given  $X$ , they become independent:

$$(16.7) \quad P(Y, Z|X) = \frac{P(X, Y, Z)}{P(X)} = \frac{P(X)P(Y|X)P(Z|X)}{P(X)} = P(Y|X)P(Z|X)$$

When its value is known,  $X$  blocks the path between  $Y$  and  $Z$ , or in other words, separates them.

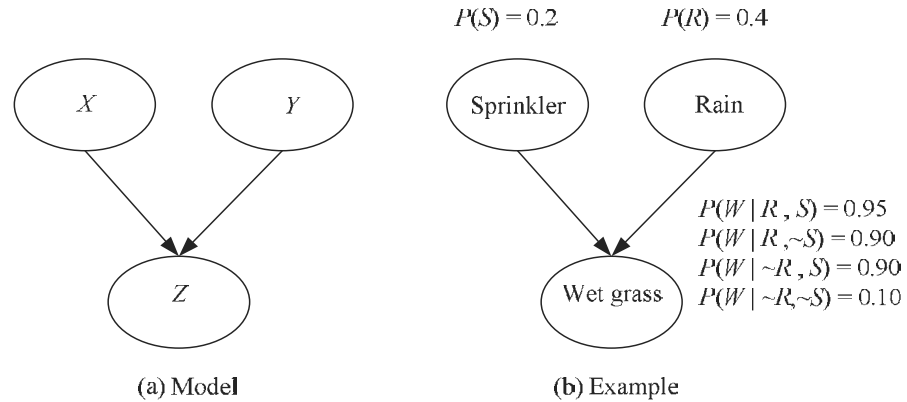
In figure 16.3b, we see an example where cloudy weather influences both rain and the use of the sprinkler, one positively and the other negatively. Knowing that it rained, for example, we can invert the dependency using Bayes' rule and infer the cause:

$$(16.8) \quad \begin{aligned} P(C|R) &= \frac{P(R|C)P(C)}{P(R)} = \frac{P(R|C)P(C)}{\sum_C P(R, C)} \\ &= \frac{P(R|C)P(C)}{P(R|C)P(C) + P(R|\sim C)P(\sim C)} = 0.89 \end{aligned}$$

Note that this value is larger than  $P(C)$ ; knowing that it rained increased the probability that the weather is cloudy.

In figure 16.3a, if  $X$  is not known, knowing  $Y$ , for example, we can infer  $X$  which we can then use to infer  $Z$ . In figure 16.3b, knowing the state of the sprinkler has an effect on the probability that it rained. If we know that the sprinkler is on,

$$(16.9) \quad P(R|S) = \sum_C P(R, C|S) = P(R|C)P(C|S) + P(R|\sim C)P(\sim C|S)$$



**Figure 16.4** Head-to-head connection. A node has two parents that are independent unless the child is given. For example, an event may have two independent causes.

$$\begin{aligned}
 &= P(R|C) \frac{P(S|C)P(C)}{P(S)} + P(R|\sim C) \frac{P(S|\sim C)P(\sim C)}{P(\sim S)} \\
 &= 0.22
 \end{aligned}$$

This is less than  $P(R) = 0.45$ ; that is, knowing that the sprinkler is on decreases the probability that it rained because sprinkler and rain happens for different states of cloudy weather. If the sprinkler is known to be off, using the same approach, we find that  $P(R|\sim S) = 0.55$ ; the probability of rain increases this time.

### Case 3: Head-to-head Connection

In a head-to-head node, there are two parents  $X$  and  $Y$  to a single node  $Z$ , as shown in figure 16.4a. The joint density is written as

$$(16.10) \quad P(X, Y, Z) = P(X)P(Y)P(Z|X, Y)$$

$X$  and  $Y$  are independent:  $P(X, Y) = P(X) \cdot P(Y)$  (exercise 2); they become dependent when  $Z$  is known. The concept of blocking or separation is different for this case: The path between  $X$  and  $Y$  is blocked, or they are separated, when  $Z$  is *not* observed; when  $Z$  (or any of its descendants) is observed, they are not blocked, separated, nor are independent.

We see for example in figure 16.4b that node  $W$  has two parents,  $R$  and  $S$ , and thus its probability is conditioned on the values of those two,  $P(W|R, S)$ .

Not knowing anything else, the probability that grass is wet is calculated by marginalizing over the joint:

$$\begin{aligned}
 P(W) &= \sum_{R,S} P(W, R, S) \\
 &= P(W|R, S)P(R, S) + P(W|\sim R, S)P(\sim R, S) \\
 &\quad + P(W|R, \sim S)P(R, \sim S) + P(W|\sim R, \sim S)P(\sim R, \sim S) \\
 &= P(W|R, S)P(R)P(S) + P(W|\sim R, S)P(\sim R)P(S) \\
 &\quad + P(W|R, \sim S)P(R)P(\sim S) + P(W|\sim R, \sim S)P(\sim R)P(\sim S) \\
 &= 0.52
 \end{aligned}$$

Now, let us say that we know that the sprinkler is on, and we check how this affects the probability. This is a causal (predictive) inference:

$$\begin{aligned}
 P(W|S) &= \sum_R P(W, R|S) \\
 &= P(W|R, S)P(R|S) + P(W|\sim R, S)P(\sim R|S) \\
 &= P(W|R, S)P(R) + P(W|\sim R, S)P(\sim R) \\
 &= 0.92
 \end{aligned}$$

We see that  $P(W|S) > P(W)$ ; knowing that the sprinkler is on, the probability of wet grass increases.

We can also calculate the probability that the sprinkler is on, given that the grass is wet. This is a diagnostic inference.

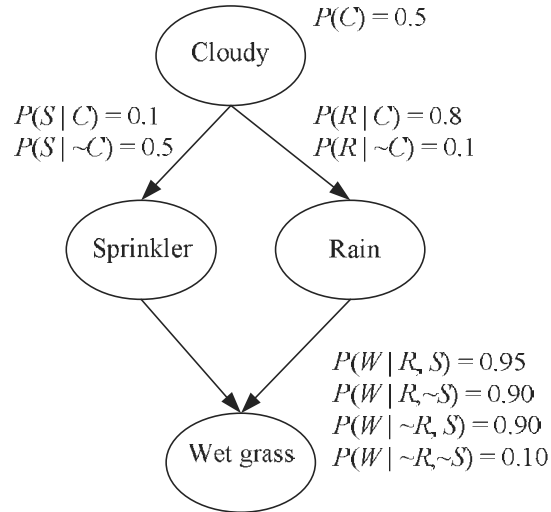
$$P(S|W) = \frac{P(W|S)P(S)}{P(W)} = 0.35$$

$P(S|W) > P(S)$ , that is, knowing that the grass is wet increased the probability of having the sprinkler on. Now let us assume that it rained. Then we have

$$\begin{aligned}
 P(S|R, W) &= \frac{P(W|R, S)P(S|R)}{P(W|R)} = \frac{P(W|R, S)P(S)}{P(W|R)} \\
 &= 0.21
 \end{aligned}$$

EXPLAINING AWAY

which is less than  $P(S|W)$ . This is called *explaining away*; given that we know it rained, the probability of sprinkler causing the wet grass decreases. Knowing that the grass is wet, rain and sprinkler become dependent. Similarly,  $P(S|\sim R, W) > P(S|W)$ . We see the same behavior when we compare  $P(R|W)$  and  $P(R|W, S)$  (exercise 3).



**Figure 16.5** Larger graphs are formed by combining simpler subgraphs over which information is propagated using the implied conditional independencies.

We can construct larger graphs by combining such subgraphs. For example, in figure 16.5 where we combine the two subgraphs, we can, for example, calculate the probability of having wet grass if it is cloudy:

$$\begin{aligned}
 P(W|C) &= \sum_{R,S} P(W,R,S|C) \\
 &= P(W,R,S|C) + P(W,\sim R,S|C) \\
 &\quad + P(W,R,\sim S|C) + P(W,\sim R,\sim S|C) \\
 &= P(W|R,S,C)P(R,S|C) \\
 &\quad + P(W|\sim R,S,C)P(\sim R,S|C) \\
 &\quad + P(W|R,\sim S,C)P(R,\sim S|C) \\
 &\quad + P(W|\sim R,\sim S,C)P(\sim R,\sim S|C) \\
 &= P(W|R,S)P(R|C)P(S|C) \\
 &\quad + P(W|\sim R,S)P(\sim R|C)P(S|C) \\
 &\quad + P(W|R,\sim S)P(R|C)P(\sim S|C) \\
 &\quad + P(W|\sim R,\sim S)P(\sim R|C)P(\sim S|C)
 \end{aligned}$$

where we have used that  $P(W|R, S, C) = P(W|R, S)$ ; given  $R$  and  $S$ ,  $W$  is independent of  $C$ :  $R$  and  $S$  between them block the path between  $W$  and  $C$ . Similarly,  $P(R, S|C) = P(R|C)P(S|C)$ ; given  $C$ ,  $R$  and  $S$  are independent. We see the advantage of Bayesian networks here, which explicitly encode independencies and allow breaking down inference into calculation over small groups of variables that are propagated from evidence nodes to query nodes.

We can calculate  $P(C|W)$  and have a diagnostic inference:

$$P(C|W) = \frac{P(W|C)P(C)}{P(W)}$$

The graphical representation is visual and helps understanding. The network represents conditional independence statements and allows us to break down the problem of representing the joint distribution of many variables into *local* structures; this eases both analysis and computation. Figure 16.5 represents a joint density of four binary variables that would normally require fifteen values ( $2^4 - 1$ ) to be stored, whereas here there are only nine. If each node has a small number of parents, the complexity decreases from exponential to linear (in the number of nodes). As we have seen earlier, inference is also easier as the joint density is broken down into conditional densities of smaller groups of variables:

$$(16.11) \quad P(C, S, R, W) = P(C)P(S|C)P(R|C)P(W|S, R)$$

In the general case, when we have variables  $X_1, \dots, X_d$ , we write

$$(16.12) \quad P(X_1, \dots, X_d) = \prod_{i=1}^d P(X_i | \text{parents}(X_i))$$

Then given any subset of  $X_i$ , namely, setting them to certain values due to evidence, we can calculate the probability distribution of some other subset of  $X_i$  by marginalizing over the joint. This is costly because it requires calculating an exponential number of joint probability combinations, even though each of them can be simplified as in equation 16.11. Note, however, that given the same evidence, for different  $X_i$ , we may be using the same intermediate values (products of conditional probabilities and sums for marginalization), and in section 16.5, we will discuss the belief propagation algorithm to do inference cheaply by doing the local intermediate calculations once which we can use multiple times for different query nodes.

Though in this example we use binary variables, it is straightforward to generalize for cases where the variables are discrete with any number of possible values (with  $m$  possible values and  $k$  parents, a table of size  $m^k$  is needed for the conditional probabilities), or they can be continuous (parameterized, e.g.,  $p(Y|x) \sim \mathcal{N}(\mu(x|\theta), \sigma^2)$ ; see section 16.3.3).

One major advantage to using a Bayesian network is that we do not need to designate explicitly certain variables as input and certain others as output. The value of any set of variables can be established through evidence and the probabilities of any other set of variables can be inferred, and the difference between unsupervised and supervised learning becomes blurry. From this perspective, a graphical model can be thought of as a “probabilistic database” (Jordan 2009), a machine that can answer queries regarding the values of random variables.

#### HIDDEN VARIABLES

In a problem, there may also be *hidden variables* whose values are never known through evidence. The advantage of using hidden variables is that the dependency structure can be more easily defined. For example, in basket analysis when we want to find the dependencies among items sold, let us say we know that there is a dependency among “baby food,” “diapers,” and “milk” in that a customer buying one of these is very much likely to buy the other two. Instead of putting (noncausal) arcs among these three, we may designate a hidden node “baby at home” as the hidden cause of the consumption of these three items. When there are hidden nodes, their values are estimated given the values of observed nodes and filled in.

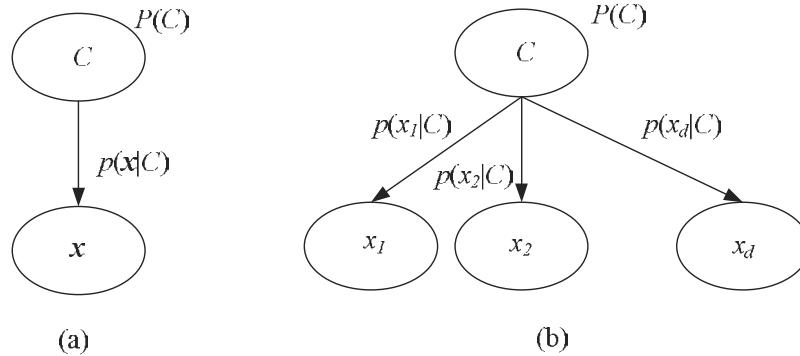
#### CAUSALITY

It should be stressed at this point that a link from a node  $X$  does not, and need not, always imply a *causality*. It only implies a *direct influence* of  $X$  over  $Y$  in the sense that the probability of  $Y$  is conditioned on the value of  $X$ , and two nodes may have a link between them even if there is no direct cause. It is preferable to have the causal relations in constructing the network by providing an explanation of how the data is generated (Pearl 2000) but such causes may not always be accessible.

## 16.3 Example Graphical Models

### 16.3.1 Naive Bayes' Classifier

For the case of classification, the corresponding graphical model is shown in figure 16.6a, with  $\mathbf{x}$  as the input and  $C$  a multinomial variable taking



**Figure 16.6** (a) Graphical model for classification. (b) Naive Bayes' classifier assumes independent inputs.

one of  $K$  states for the class code. Bayes' rule allows a diagnosis, as in the rain and wet grass case we saw in figure 16.1:

$$P(C|\mathbf{x}) = \frac{P(C)p(\mathbf{x}|C)}{P(\mathbf{x})}$$

NAIVE BAYES'  
CLASSIFIER

If the inputs are independent, we have the graph shown in figure 16.6b, which is called the *naive Bayes' classifier*, because it ignores possible dependencies, namely, correlations, among the inputs and reduces a multivariate problem to a group of univariate problems:

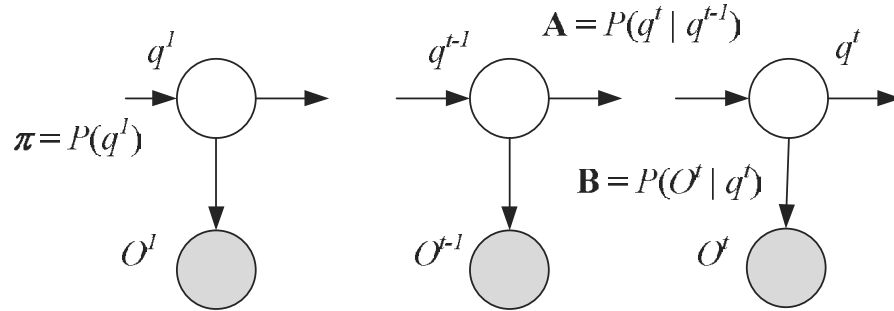
$$p(\mathbf{x}|C) = \prod_{j=1}^d p(x_j|C)$$

We have discussed classification for this case in sections 5.5 and 5.7 for numeric and discrete  $\mathbf{x}$ , respectively.

Clustering is also similar except that the multinomial class indicator variable  $C$  is observed in classification, but the similar variable,  $Z$ , cluster indicator, is not observed. The E-step of the Expectation Maximization algorithm (section 7.4) uses Bayes' rule to invert the arc and estimates the cluster indicator given the input.

GENERATIVE MODEL

Figure 16.6a is a *generative model* of the process that creates the data. It is as if we first pick a class  $C$  at random by sampling from  $P(C)$ , and then having fixed  $C$ , we pick an  $\mathbf{x}$  by sampling from  $p(\mathbf{x}|C)$ . Thinking of data as sampled from a causal generative model that can be visualized as a graph can ease understanding and also inference in many domains.



**Figure 16.7** Hidden Markov model can be drawn as a graphical model where  $q^t$  are the hidden states and shaded  $O^t$  are observed.

For example, in text categorization, generating a text may be thought of as the process where an author decides to write a document on a certain topic and then chooses the set of words accordingly. In bioinformatics, one area among many where a graphical approach used is the modeling of a *phylogenetic tree*; namely, a directed graph whose leaves are the current species, nonterminal nodes are past ancestors that split into multiple species during a speciation event, and the conditional probabilities depend on the evolutionary distance between a species and its ancestor (Jordan 2004).

PHYLOGENETIC TREE

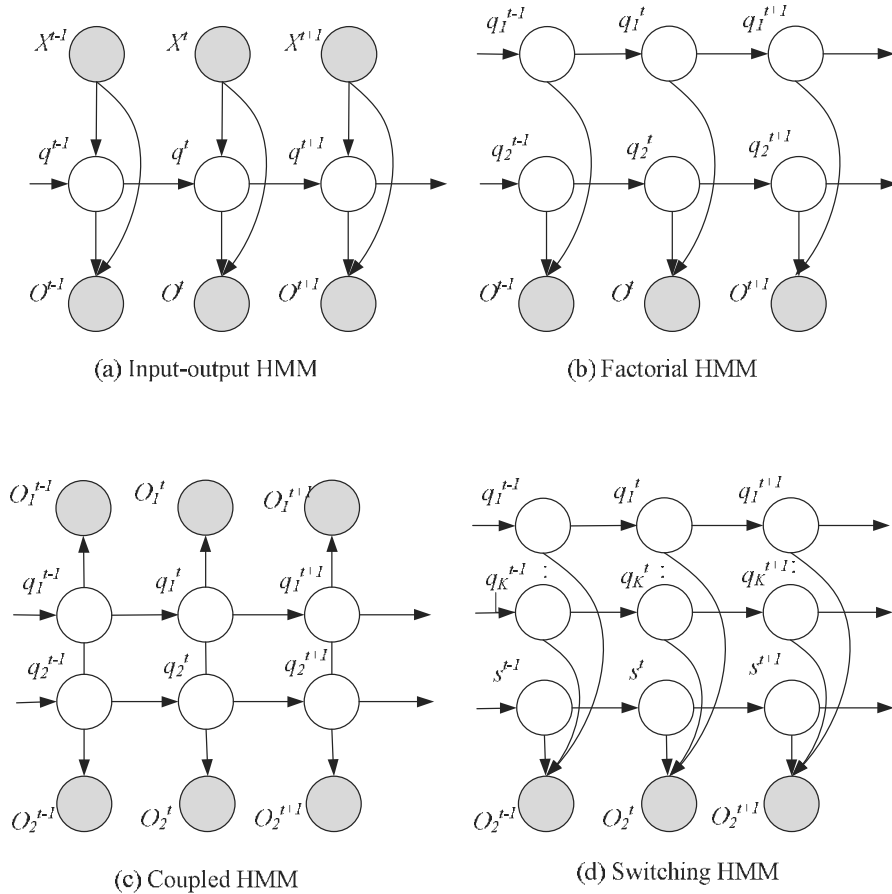
### 16.3.2 Hidden Markov Model

HIDDEN MARKOV MODEL

*Hidden Markov models* (HMM), which we previously discussed in chapter 15, are an example of case 1 where three successive states  $q_{t-2}, q_{t-1}, q_t$  correspond to three states on a chain in a first-order Markov model. The state at time  $t$ ,  $q_t$ , depends only on the state at time  $t - 1$ ,  $q_{t-1}$ , and given  $q_{t-1}$ ,  $q_t$  is independent of  $q_{t-2}$

$$P(q_t | q_{t-1}, q_{t-2}) = P(q_t | q_{t-1})$$

as given by the state transition probability matrix **A** (see figure 16.7). Each hidden variable generates a discrete observation that is observed, as given by the observation probability matrix **B**. The forward-backward procedure of hidden Markov models is a special case of belief propagation that we will discuss shortly.



**Figure 16.8** Different types of HMM model different assumptions about the way the observed data (shown shaded) is generated from Markov sequences of latent variables.

INPUT-OUTPUT HMM

Different HMM types can be shown as different graphical models. In figure 16.8a, an *input-output HMM* is shown (see section 15.9) where there are two separate observed input-output sequences and there is also a sequence of hidden states. The output observation depends both on the state and also on the input; one can think of this as a **B** matrix whose elements are not scalars but parametrized functions of the input. This may similarly be seen as a mixture of expert architecture (section 12.8)

whose gating output (hidden state) depends also on the gating value at the previous time step.

FACTORIAL HMM

Another HMM type that can be easily visualized is a *factorial HMM*, where there are multiple separate hidden sequences that interact to generate a single observation sequence. An example is a *pedigree* which displays the parent-child relationship (Jordan 2004); figure 16.8b models *meiosis* where the two sequences correspond to the chromosomes of the father and the mother (which are independent), and at each locus (gene), the offspring receives one allele from the father or the other allele from the mother.

PEDIGREE

COUPLED HMM

A *coupled HMM*, shown in figure 16.8c, models two parallel but related hidden sequences that generate two parallel observation sequences. For example, in speech recognition, we may have one observed acoustic sequence of uttered words and one observed visual sequence of lip images, each having its hidden states where the two are dependent.

SWITCHING HMM

In a *switching HMM*, shown in figure 16.8d, there are  $K$  parallel independent hidden state sequences and the state variable  $S$  at any one time picks one of them and the chosen one generates the output. That is, we switch between state sequences as we go along.

LINEAR DYNAMICAL  
SYSTEM  
KALMAN FILTER

In HMM proper, though the observation may be continuous, state is discrete; in a *linear dynamical system*, also known as the *Kalman filter*, both the state and the observations are continuous. In the basic case, state at time  $t$  is a linear function of state at  $t - 1$  with additive zero-mean Gaussian noise, and, at each state, the observation is another linear function of the state with additive zero-mean Gaussian noise. The two linear mappings and the covariances of the two noise sources make up the parameters. All HMM variants we discussed earlier can similarly be generalized to use continuous states.

By suitably modifying the graphical model, one can adapt the architecture to the characteristics of the process that generates the data. This process of matching the model to the data is a model selection procedure to best trade off bias and variance. The disadvantage is that exact inference may no longer be possible on such extended HMMs, and one would need approximation or sampling methods (Ghahramani 2001; Jordan 2009).

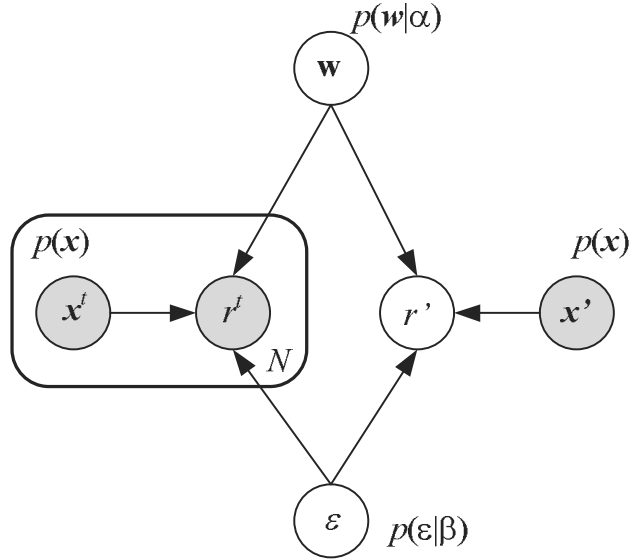


Figure 16.9 Bayesian network for linear regression.

### 16.3.3 Linear Regression

Linear regression can be visualized as a graphical model, as shown in figure 16.9. Input  $\mathbf{x}^t$  is drawn from a prior  $p(\mathbf{x})$  and the dependent variable  $r^t$  depend on the input  $\mathbf{x}$ , weights  $\mathbf{w}$  (drawn from a prior parameterized by  $\alpha$ , i.e.,  $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I})$ ), and noise  $\epsilon$  (parameterized by  $\beta$ , i.e.,  $p(\epsilon) \sim \mathcal{N}(0, \beta^{-1})$ ):

$$(16.13) \quad p(r^t | \mathbf{x}^t, \mathbf{w}) \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}^t, \beta^{-1})$$

There are  $N$  such pairs in the training set, which is shown by the rectangular *plate* in the figure. Given a new input  $\mathbf{x}'$ , the aim is to estimate  $r'$ , which will be  $E[r' | \mathbf{x}', \mathbf{w}]$ .

The weights  $\mathbf{w}$  are not given but they can be estimated using the training set of  $[\mathbf{X}, \mathbf{r}]$ . Just as in equation 16.9, where  $C$  is the cause of  $R$  and  $S$ , where we used

$$P(R|S) = \sum_C P(R, C|S) = P(R|C)P(C|S) + P(R|\sim C)P(\sim C|S)$$

filling in  $C$  using  $S$ , which we in turn used to estimate  $R$ . Here, we write

$$\begin{aligned}
 p(r'|\mathbf{x}', \mathbf{r}, \mathbf{X}) &= \int p(r'|\mathbf{x}', \mathbf{w})p(\mathbf{w}|\mathbf{X}, \mathbf{r})d\mathbf{w} \\
 &= \int p(r'|\mathbf{x}', \mathbf{w})\frac{p(\mathbf{r}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{r})}d\mathbf{w} \\
 (16.14) \qquad &\propto \int p(r'|\mathbf{x}', \mathbf{w})\prod_t p(r^t|\mathbf{x}^t, \mathbf{w})d\mathbf{w}
 \end{aligned}$$

where the second line is due to Bayes' rule and the third line is due to the independence of instances in the training set.

## 16.4 d-Separation

D-SEPARATION

BAYES' BALL

We now generalize the concept of blocking and separation under the name of *d-separation*, and we define it in a way so that for arbitrary subsets of nodes  $A$ ,  $B$ , and  $C$ , we can check if  $A$  and  $B$  are independent given  $C$ . Jordan (2009) visualizes this as a ball bouncing over the graph and calls this the *Bayes' ball*. We set the nodes in  $C$  to their values, place a ball at each node in  $A$ , let the balls move around according to a set of rules, and check whether a ball reaches any node in  $B$ . If this is the case, they are dependent; otherwise, they are independent.

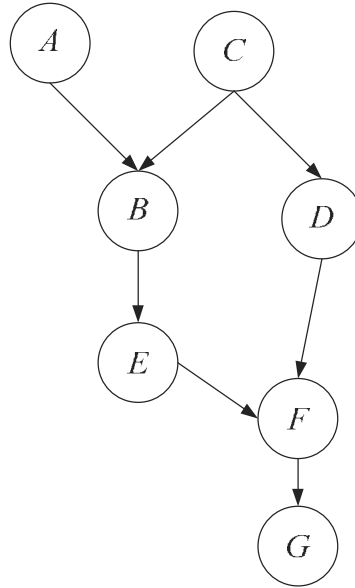
To check whether  $A$  and  $B$  are d-separated given  $C$ , we consider all possible paths between any node in  $A$  and any node in  $B$ . Any such path is *blocked* if

- (a) the directions of the edges on the path either meet head-to-tail (case 1) or tail-to-tail (case 2) and the node is in  $C$ , or
- (b) the directions of the edges on the path meet head-to-head (case 3) and neither that node nor any of its descendant is in  $C$ .

If all paths are blocked, we say that  $A$  and  $B$  are d-separated, that is, independent, given  $C$ ; otherwise, they are dependent. Examples are given in figure 16.10.

## 16.5 Belief Propagation

Having discussed some inference examples by hand, we now are interested in an algorithm that can answer queries such as  $P(X|E)$  where  $X$



**Figure 16.10** Examples of d-separation. The path  $BCDF$  is blocked given  $C$  because  $C$  is a tail-to-tail node.  $BEFG$  is blocked by  $F$  because  $F$  is a head-to-tail node.  $BEFD$  is blocked unless  $F$  (or  $G$ ) is given.

is any *query node* in the graph and  $E$  is any subset of *evidence nodes* whose values are set to certain value. Following Pearl (1988), we start with the simplest case of chains and gradually move on to more complex graphs. Our aim is to find the graph operation counterparts of probabilistic procedures such as Bayes' rule or marginalization, so that the task of inference can be mapped to general purpose graph algorithms.

### 16.5.1 Chains

A *chain* is a sequence of head-to-tail nodes with one *root* node without any parent; all other nodes have exactly one parent node, and all nodes except the very last, *leaf*, have a single child. If evidence is in the ancestors of  $X$ , we can just do a diagnostic inference and propagate evidence down the chain; if evidence is in the descendants of  $X$ , we can do a causal inference and propagate upward using Bayes' rule. Let us see the general case where we have evidence in both directions, up the chain  $E^+$  and

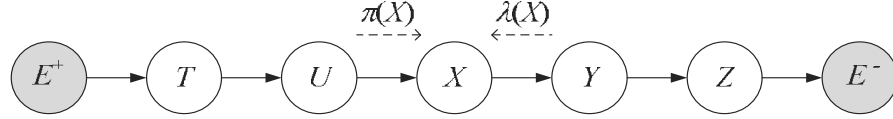


Figure 16.11 Inference along a chain.

down the chain  $E^-$  (see figure 16.11). Note that any evidence node separates  $X$  from the nodes on the chain on the other side of the evidence and their values do not affect  $p(X)$ ; this is true in both directions.

We consider each node as a processor that receives messages from its neighbors and pass it along after some local calculation. Each node  $X$  locally calculates and stores two values:  $\lambda(X) \equiv P(E^-|X)$  is the propagated  $E^-$  that  $X$  receives from its child and forwards to its parent, and  $\pi(X) \equiv P(X|E^+)$  is the propagated  $E^+$  that  $X$  receives from its parent and passes on to its child.

$$\begin{aligned}
 P(X|E) &= \frac{P(E|X)P(X)}{P(E)} = \frac{P(E^+, E^-|X)P(X)}{P(E)} \\
 &= \frac{P(E^+|X)P(E^-|X)P(X)}{P(E)} \\
 &= \frac{P(X|E^+)P(E^+)P(E^-|X)P(X)}{P(X)P(E)} \\
 (16.15) \quad &= \alpha P(X|E^+)P(E^-|X) = \alpha \pi(X)\lambda(X)
 \end{aligned}$$

for some normalizing constant  $\alpha$ , not dependent on the value of  $X$ . The second line is there because  $E^+$  and  $E^-$  are independent given  $X$ , and the third line is due to Bayes' rule.

If a node  $E$  is instantiated to a certain value  $\tilde{e}$ ,  $\lambda(\tilde{e}) \equiv 1$  and  $\lambda(e) \equiv 0$ , for  $e \neq \tilde{e}$ . The leaf node  $X$  that is not instantiated has its  $\lambda(x) \equiv 1$ , for all  $x$  values. The root node  $X$  that is not instantiated takes the prior probabilities as  $\pi$  values:  $\pi(x) \equiv P(x)$ ,  $\forall x$ .

Given these initial conditions, we can devise recursive formulas to propagate evidence along the chain.

For the  $\pi$ -messages, we have

$$\begin{aligned}
 \pi(X) &\equiv P(X|E^+) = \sum_U P(X|U, E^+)P(U|E^+) \\
 (16.16) \quad &= \sum_U P(X|U)P(U|E^+) = \sum_U P(X|U)\pi(U)
 \end{aligned}$$

where the second line follows from the fact that  $U$  blocks the path between  $X$  and  $E^+$ .

For the  $\lambda$ -messages, we have

$$\begin{aligned}
 \lambda(X) &\equiv P(E^-|X) = \sum_Y P(E^-|X, Y)P(Y|X) \\
 (16.17) \quad &= \sum_Y P(E^-|Y)P(Y|X) = \sum_U P(Y|X)\lambda(Y)
 \end{aligned}$$

where the second line follows from the fact that  $Y$  blocks the path between  $X$  and  $E^-$ .

When the evidence nodes are set to a value, they initiate traffic and nodes continue updating until there is convergence. Pearl (1988) views this as a parallel machine where each node is implemented by a processor that works in parallel with others and exchanges information through  $\lambda$ - and  $\pi$ -messages with its parent and child.

### 16.5.2 Trees

Chains are restrictive because each node can have only a single parent and a single child, that is, a single cause and a single symptom. In a *tree*, each node may have several children but all nodes, except the single root, have exactly one parent. The same belief propagation also applies here with the difference from chains being that a node receives different  $\lambda$ -messages from its children,  $\lambda_Y(X)$  denoting the message  $X$  receives from its child  $Y$ , and sends different  $\pi$ -messages to its children,  $\pi_Y(X)$  denoting the message  $X$  sends to its child  $Y$ .

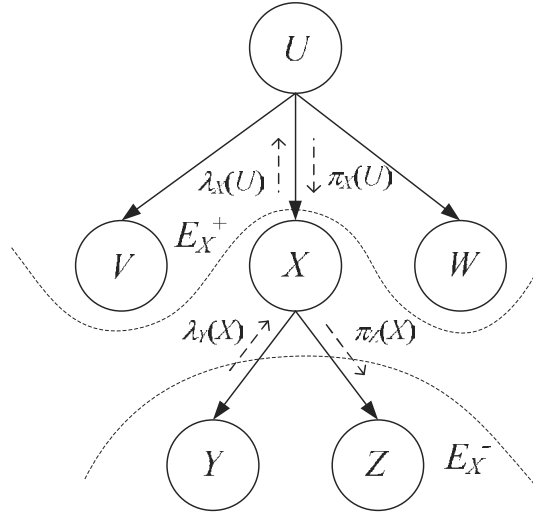
Again, we divide possible evidence to two parts,  $E^-$  are nodes that are in the subtree rooted at the query node  $X$ , and  $E^+$  are evidence nodes elsewhere (see figure 16.12). Note that this second need not be an ancestor of  $X$  but may also be in a subtree rooted at a sibling of  $X$ . The important point is that again  $X$  separates  $E^+$  and  $E^-$  so that we can write  $P(E^+, E^-|X) = P(E^+|X)P(E^-|X)$ , and hence have

$$P(X|E) = \alpha \pi(X) \lambda(X)$$

where again  $\alpha$  is a normalizing constant.

$\lambda(X)$  is the evidence in the subtree rooted at  $X$ , and if  $X$  has two children  $Y$  and  $Z$ , as shown in figure 16.12, it can be calculated as

$$\begin{aligned}
 \lambda(X) &\equiv P(E_X^-|X) = P(E_Y^-, E_Z^-|X) \\
 (16.18) \quad &= P(E_Y^-|X)P(E_Z^-|X) = \lambda_Y(X)\lambda_Z(X)
 \end{aligned}$$



**Figure 16.12** In a tree, a node may have several children but a single parent.

In the general case if  $X$  has  $m$  children,  $Y_j, j = 1, \dots, m$ , then we multiply all their  $\lambda$  values:

$$(16.19) \quad \lambda(X) = \prod_{j=1}^m \lambda_{Y_j}(X)$$

Once  $X$  accumulates  $\lambda$  evidence from its children's  $\lambda$ -messages, it propagates it up to its parent:

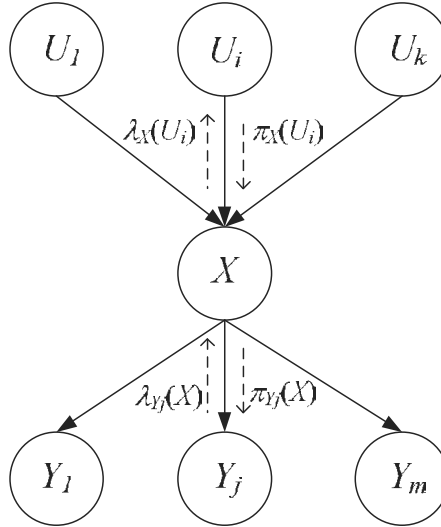
$$(16.20) \quad \lambda_X(U) = \sum_X \lambda(X) P(X|U)$$

Similarly and in the other direction,  $\pi(X)$  is the evidence elsewhere that is accumulated in  $P(U|E^+)$  and passed on to  $X$  as a  $\pi$ -message:

$$(16.21) \quad \pi(X) \equiv P(X|E_X^+) = \sum_U P(X|U) P(U|E_X^+) = \sum_U P(X|U) \pi_X(U)$$

This calculated  $\pi$  value is then propagated down to  $X$ 's children. Note that what  $Y$  receives from  $X$  is what  $X$  receives from its parent  $U$  and also from its other child  $Z$ ; together they make up  $E_Y^+$  (see figure 16.12):

$$\pi_Y(X) \equiv P(X|E_Y^+) = P(X|E_X^+, E_Z^-)$$



**Figure 16.13** In a polytree, a node may have several children and several parents, but the graph is singly connected; that is, there is a single chain between  $U_i$  and  $Y_j$  passing through  $X$ .

$$\begin{aligned}
 &= \frac{P(E_{\bar{Z}}|X, E_X^+)P(X|E_X^+)}{P(E_{\bar{Z}})} = \frac{P(E_{\bar{Z}}|X)P(X|E_X^+)}{P(E_{\bar{Z}})} \\
 (16.22) \quad &= \alpha \lambda_Z(X) \pi(X)
 \end{aligned}$$

Again, if  $Y$  has not one sibling  $Z$  but multiple, we need to take a product over all their  $\lambda$  values:

$$(16.23) \quad \pi_{Y_j}(X) = \alpha \prod_{s \neq j} \lambda_{Y_s}(X) \pi(X)$$

### 16.5.3 Polytrees

POLYTREE In a tree, a node has a single parent, that is, a single cause. In a *polytree*, a node may have multiple parents, but we require that the graph be singly connected, which means that there is a single chain between any two nodes. If we remove  $X$ , the graph will split into two components. This is necessary so that we can continue splitting  $E_X$  into  $E_X^+$  and  $E_X^-$ , which are independent given  $X$  (see figure 16.13).

If  $X$  has multiple parents  $U_i, i = 1, \dots, k$ , it receives  $\pi$ -messages from

all of them,  $\pi_X(U_i)$ , which it combines as follows:

$$\begin{aligned}
 \pi(X) &\equiv P(X|E_X^+) = P(X, E_{U_1X}^+, E_{U_2X}^+, \dots, E_{U_kX}^+) \\
 &= \sum_{U_1} \sum_{U_2} \cdots \sum_{U_k} P(X|U_1, U_2, \dots, U_k) P(U_1|E_{U_1X}^+) \cdots P(U_k|E_{U_kX}^+) \\
 (16.24) \quad &= \sum_{U_1} \sum_{U_2} \cdots \sum_{U_k} P(X|U_1, U_2, \dots, U_k) \prod_{i=1}^k \pi_X(U_i)
 \end{aligned}$$

and passes it on to its several children  $Y_j, j = 1, \dots, m$ :

$$(16.25) \quad \pi_{Y_j}(X) = \alpha \prod_{s \neq j} \lambda_{Y_s}(X) \pi(X)$$

In this case when  $X$  has multiple parents, a  $\lambda$ -message  $X$  passes on to one of its parents  $U_i$  combines not only the evidence  $X$  receives from its children but also the  $\pi$ -messages  $X$  receives from its other parents  $U_r, r \neq i$ ; they together make up  $E_{U_iX}^-$ :

$$\begin{aligned}
 \lambda_X(U_i) &\equiv P(E_{U_iX}^-|X) \\
 &= \sum_X \sum_{U_{r \neq i}} P(E_X^-, E_{U_{r \neq i}X}^+, X, U_{r \neq i}|U_i) \\
 &= \sum_X \sum_{U_{r \neq i}} P(E_X^-, E_{U_{r \neq i}X}^+|X, U_{r \neq i}, U_i) P(X, U_{r \neq i}|U_i) \\
 &= \sum_X \sum_{U_{r \neq i}} P(E_X^-|X) P(E_{U_{r \neq i}X}^+|U_{r \neq i}) P(X|U_{r \neq i}, U_i) P(U_{r \neq i}|U_i) \\
 &= \sum_X \sum_{U_{r \neq i}} P(E_X^-|X) \frac{P(U_{r \neq i}|E_{U_{r \neq i}X}^+) P(E_{U_{r \neq i}X}^+)}{P(U_{r \neq i})} P(X|U_{r \neq i}, U_i) P(U_{r \neq i}|U_i) \\
 &= \beta \sum_X \sum_{U_{r \neq i}} P(E_X^-|X) P(U_{r \neq i}|E_{U_{r \neq i}X}^+) P(X|U_{r \neq i}, U_i) \\
 &= \beta \sum_X \sum_{U_{r \neq i}} \lambda(X) \prod_{r \neq i} \pi_X(U_r) P(X|U_1, \dots, U_k) \\
 (16.26) \quad &= \beta \sum_X \lambda(X) \sum_{U_{r \neq i}} P(X|U_1, \dots, U_k) \prod_{r \neq i} \pi_X(U_r)
 \end{aligned}$$

As in a tree, to find its overall  $\lambda$ , the parent multiplies the  $\lambda$ -messages it receives from its children:

$$(16.27) \quad \lambda(X) = \prod_{j=1}^m \lambda_{Y_j}(X)$$

In this case of multiple parents, we need to store and manipulate the conditional probability given all the parents,  $p(X|U_1, \dots, U_k)$ , which is costly for large  $k$ . Approaches have been proposed to decrease the complexity from exponential in  $k$  to linear. For example, in a *noisy OR gate*, any of the parents is sufficient to cause the event and the likelihood does not decrease when multiple parent events occur. If the probability that  $X$  happens when only cause  $U_i$  happens is  $1 - q_i$

$$(16.28) \quad P(X|U_i, \sim U_{p \neq j}) = 1 - q_i$$

the probability that  $X$  happens when a subset  $T$  of them occur is calculated as

$$(16.29) \quad P(X|T) = 1 - \prod_{u_i \in T} q_i$$

For example, let us say wet grass has two causes, rain and a sprinkler, with  $q_R = q_S = 0.1$ ; that is, both singly have a 90 percent probability of causing wet grass. Then,  $P(W|R, \sim S) = 0.9$  and  $P(W|R, S) = 0.99$ .

Another possibility is to write the conditional probability as some function given a set of parameters, for example, as a linear model

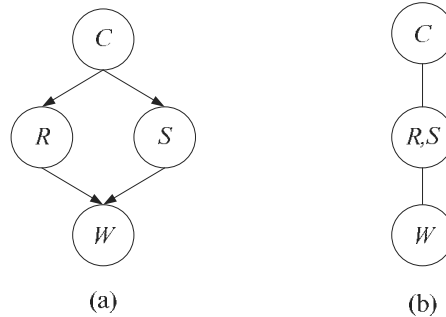
$$(16.30) \quad P(X|U_1, \dots, U_k, w_0, w_1, \dots, w_k) = \text{sigmoid} \left( \sum_{i=1}^k w_i U_i + w_0 \right)$$

where sigmoid guarantees that the output is a probability between 0 and 1. During training, we can learn the parameters  $w_i, i = 0, \dots, d$ , for example, to maximize the likelihood on a sample.

#### 16.5.4 Junction Trees

If there is a loop, that is, if there is a cycle in the underlying undirected graph—for example, if the parents of  $X$  share a common ancestor—the algorithm we discussed earlier does not work. In such a case, there is more than one path on which to propagate evidence and, for example, while evaluating the probability at  $X$ , we cannot say that  $X$  separates  $E$  into  $E_X^+$  and  $E_X^-$  as causal (upward) and diagnostic (downward) evidence; removing  $X$  does not split the graph into two. Conditioning them on  $X$  does not make them independent and the two can interact through some other path not involving  $X$ .

We can still use the same algorithm if we can convert the graph to a polytree. We define *clique nodes* that correspond to a set of original variables and connect them so that they form a tree (see figure 16.14). We



**Figure 16.14** (a) A multiply connected graph, and (b) its corresponding junction tree with nodes clustered.

JUNCTION TREE

can then run the same belief propagation algorithm with some modifications. This is the basic idea behind the *junction tree algorithm* (Lauritzen and Spiegelhalter 1988; Jensen 1996; Jordan 2009).

## 16.6 Undirected Graphs: Markov Random Fields

MARKOV RANDOM  
FIELD

Up to now, we have discussed directed graphs where the influences are unidirectional and have used Bayes' rule to invert the arcs. If the influences are symmetric, we represent them using an undirected graphical model, also known as a *Markov random field*. For example, neighboring pixels in an image tend to have the same color—that is, are correlated—and this correlation goes both ways.

Directed and undirected graphs define conditional independence differently, and, hence, there are probability distributions that are represented by a directed graph and not by an undirected graph, and vice versa (Pearl 1988).

Because there are no directions and hence no distinction between the head or the tail of an arc, the treatment of undirected graphs is simpler. For example, it is much easier to check if  $A$  and  $B$  are independent given  $C$ . We just check if after removing all nodes in  $C$ , we still have a path between a node in  $A$  and a node in  $B$ . If so, they are dependent, otherwise, if all paths between nodes in  $A$  and nodes in  $B$  pass through nodes in  $C$  such that removal of  $C$  leaves nodes of  $A$  and nodes of  $B$  in separate components, we have independence.

POTENTIAL FUNCTION

CLIQUE

In the case of an undirected graph, we do not talk about the parent or the child but about *cliques*, which are sets of nodes such that there exists a link between any two nodes in the set. A *maximal* clique has the maximum number of elements. Instead of conditional probabilities (implying a direction), in undirected graphs we have *potential functions*  $\psi_C(X_C)$  where  $X_C$  is the set of variables in clique  $C$ , and we define the joint distribution as the product of the potential functions of the maximal cliques of the graph

$$(16.31) \quad p(X) = \frac{1}{Z} \prod_C \psi_C(X_C)$$

where  $Z$  is the normalization constant to make sure that  $\sum_X p(X) = 1$ :

$$(16.32) \quad Z = \sum_X \prod_C \psi_C(X)$$

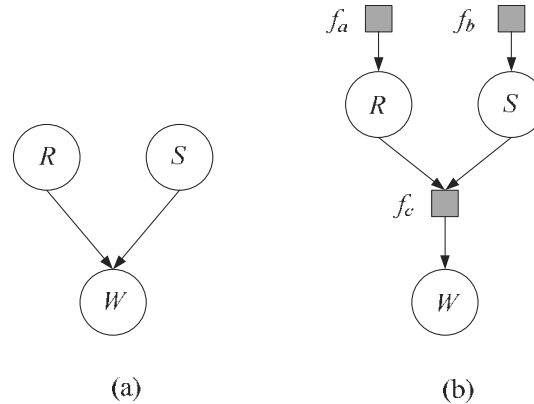
It can be shown that a directed graph is already normalized (exercise 5).

Unlike in directed graphs, the potential functions in an undirected graph do not need to have a probabilistic interpretation, and one has more freedom in defining them. In general, we can view potential functions as expressing local constraints, that is, favoring some local configurations over others. For example, in an image, we can define a pairwise potential function between neighboring pixels, which takes a higher value if their colors are similar than the case when they are different (Bishop 2006). Then, setting some of the pixels to their values given as evidence, we can estimate the values of other pixels that are not known, for example, due to occlusion.

MORALIZATION

If we have the directed graph, it is easy to redraw it as an undirected graph, simply by dropping all the directions, and if a node has a single parent, we can set the pairwise potential function simply to the conditional probability. If the node has more than one parent, however, the “explaining away” phenomenon due to the head-to-head node makes the parents dependent, and hence we should have the parents in the same clique so that the clique potential includes all the parents. This is done by connecting all the parents of a node by links so that they are completely connected among them and form a clique. This is called “marrying” the parents, and the process is called *moralization*. Incidentally, moralization is one of the steps in generating a junction tree, which is undirected.

It is straightforward to adapt the belief propagation algorithm to work on undirected graphs, and it is easier because the potential function is



**Figure 16.15** (a) A directed graph that would have a loop after moralization, and (b) its corresponding factor graph that is a tree. The three factors are  $f_a(R) \equiv P(R)$ ,  $f_b(S) \equiv P(S)$ , and  $f_c(R, S, W) \equiv P(W|R, S)$ .

FACTOR GRAPH

symmetric and we do not need to make a difference between causal and diagnostic evidence. Thus, we can do inference on undirected chains and trees. But in polytrees where a node has multiple parents and moralization necessarily creates loops, this would not work. One trick is to convert it to a *factor graph* that uses a second kind of *factor nodes* in addition to the variable nodes, and we write the joint distribution as a product of factors (Kschischang, Frey, and Loeliger 2001)

$$(16.33) \quad p(X) = \frac{1}{Z} \prod_S f_S(X_S)$$

where  $X_S$  denotes a subset of the variable nodes used by factor  $S$ . Directed graphs are a special case where factors correspond to local conditional distributions, and undirected graphs are another special case where factors are potential functions over maximal cliques. The advantage is that, as we can see in figure 16.15, the tree structure can be kept even after moralization.

SUM-PRODUCT  
ALGORITHM

It is possible to generalize the belief propagation algorithm to work on factor graphs; this is called the *sum-product algorithm* (Bishop 2006; Jordan 2009) where there is the same idea of doing local computations once and propagating them through the graph as messages. The difference now is that there are two types of messages because there are two kinds of nodes, factors and variables, and we make a distinction between their

messages. Note, however, that a factor graph is bipartite and one kind of node can have a close encounter only with the second kind.

In belief propagation, or the sum-product algorithm, the aim is to find the probability of a set of nodes  $X$  given that another set of evidence nodes  $E$  are clamped to a certain value, that is,  $P(X|E)$ . In some applications, we may be interested in finding the setting of all  $X$  that maximizes the full joint probability distribution  $p(X)$ . For example, in the undirected case where potential functions code locally consistent configurations, such an approach would propagate local constraints over the whole graph and find a solution that maximizes global consistency. In a graph where nodes correspond to pixels and pairwise potential functions favor correlation, this approach would implement noise removal (Bishop 2006). The algorithm for this, named the *max-product algorithm* (Bishop 2006; Jordan 2009) is the same as the sum-product algorithm where we take the maximum (choose the most likely value) rather than the sum (marginalize). This is analogous to the difference between the forward-backward procedure and the Viterbi algorithm in hidden Markov models that we discussed in chapter 15.

MAX-PRODUCT  
ALGORITHM

Note that the nodes need not correspond to low-level concepts like pixels; in a vision application, for instance, we may have nodes for corners of different types or lines of different orientations with potential functions checking for compatibility, so as to see if they can be part of the same interpretation—remember the Necker cube, for example—so that overall consistent solutions emerge after the consolidation of local evidences.

The complexity of the inference algorithms on polytrees or junction trees is determined by the maximum number of parents or the size of the largest clique, and when this is large, exact inference may be infeasible. In such a case, one needs to use an approximation or a sampling algorithm (Jordan 1999; Bishop 2006; Jordan 2009).

## 16.7 Learning the Structure of a Graphical Model

As in any approach, learning a graphical model has two parts. The first is the learning of parameters given a structure; this is relatively easier (Buntine 1996), and, in graphical models, conditional probability tables or their parameterizations (as in equation 16.30) can be trained to maximize the likelihood, or by using a Bayesian approach if suitable priors are known (chapter 14).

The second, more difficult, and interesting part is to learn the graph structure (Cowell et al. 1999). This is basically a model selection problem, and just like the incremental approaches for learning the structure of a multilayer perceptron (section 11.9), we can see this as a search in the space of all possible graphs. One can, for example, consider operators that can add/remove arcs and/or hidden nodes and then do a search evaluating the improvement at each step (using parameter learning at each intermediate iteration). Note, however, that to check for overfitting, one should regularize properly, corresponding to a Bayesian approach with a prior that favors simpler graphs (Neapolitan 2004). However, because the state space is large, it is most helpful if there is a human expert who can manually define causal relationships among variables and creates subgraphs of small groups of variables.

## 16.8 Influence Diagrams

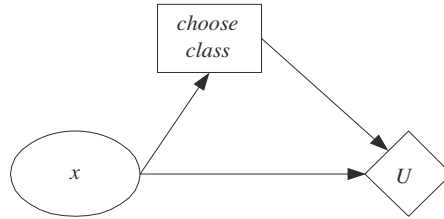
### INFLUENCE DIAGRAMS

Just as in chapter 3, we generalized from probabilities to actions with risks, *influence diagrams* are graphical models that allow the generalization of graphical models to include decisions and utilities. An influence diagram contains *chance nodes* representing random variables that we use in graphical models (see figure 16.16). It also has decision nodes and a utility node. A *decision node* represents a choice of actions. A *utility node* is where the utility is calculated. Decisions may be based on chance nodes and may affect other chance nodes and the utility node.

Inference on an influence diagram is an extension to belief propagation on a graphical model. Given evidence on some of the chance nodes, this evidence is propagated, and for each possible decision, the utility is calculated and the decision having the highest utility is chosen. The influence diagram for classification of a given input is shown in figure 16.16. Given the input, the decision node decides on a class, and for each choice we incur a certain utility (risk).

## 16.9 Notes

Graphical models have two advantages. One is that we can visualize the interaction of variables and have a better understanding of the process, for example, by using a causal generative model. The second is that by finding graph operations that correspond to basic probabilistic proce-



**Figure 16.16** Influence diagram corresponding to classification. Depending on input  $x$ , a class is chosen that incurs a certain utility (risk).

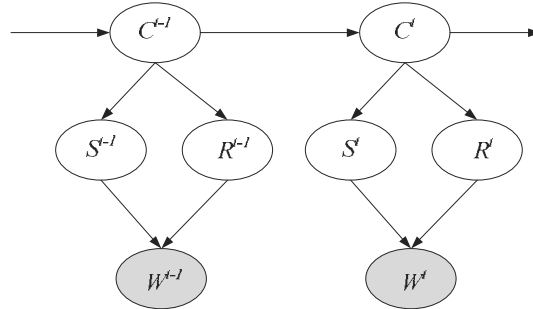
dures such as Bayes' rule or marginalization, the task of inference can be mapped to general-purpose graph algorithms that can be efficiently represented and implemented.

The idea of visual representation of variables and dependencies between them as a graph, and the related factorization of a complicated global function of many variables as a product of local functions involving a small subset of the variables for each, seems to be used in different domains in decision making, coding, and signal processing; Kschischang, Frey, and Loeliger (2001) give a review.

The complexity of the inference algorithms on polytrees or junction trees is determined by the maximum number of parents or the size of the largest clique, and when this is large exact inference may be infeasible. In such a case, one needs to use an approximation or a sampling algorithm. Variational approximations and Markov chain Monte Carlo (MCMC) algorithms are discussed in Jordan et al. 1999, MacKay 2003, Andrieu et al. 2003, Bishop 2006, and Jordan 2009.

Graphical models are especially suited to represent Bayesian approaches where in addition to nodes for variables, we also have nodes for hidden parameters that influence the observed variables. We may also introduce a hierarchy where we have nodes for the hyperparameters—that is, second-level parameters for the priors of the first-level parameters, and so on.

Hidden Markov models is one type of graphical model, and actually any graphical model can be extended in time by unfolding it in time and adding dependencies between successive copies. Such *dynamic graphical models* find application in areas where there is also a temporal dimension—speech recognition, for example. In fact, a hidden Markov model is not-



**Figure 16.17** A dynamic version where we have a chain of graphs to show dependency in weather in consecutive days.

ing but a sequence of clustering problems where the cluster index at time  $t$  is dependent not only on observation at time  $t$  but also on the index at time  $t - 1$ , and Baum-Welch algorithm is Expectation-Maximization extended to also include this dependency in time. In section 6.4, we discussed factor analysis where a small number of hidden factors generate the observation; similarly, a linear dynamical system may be viewed as a sequence of such factor analysis models where the current factors also depend on the previous factors.

This dynamic dependency may be added when needed. For example, figure 16.5 models the cause of wet grass for a particular day; if we believe that yesterday's weather has an influence on today's weather (and we should—it tends to be cloudy on successive days, then sunny for a number of days, and so on), we can have the dynamic graphical model shown in figure 16.17 where we model this dependency.

The general graphical model formalism allows us to go beyond the power of HMM proper and lead to improved performances, for example, in speech recognition (Zweig 2003; Bilmes and Bartels 2005). Graphical models are also used in computer vision—for example, in information retrieval (Barnard et al. 2003) and scene analysis (Sudderth et al. 2008). A review of the use of graphical models in bioinformatics (and related software) is given in Donkers and Tuyls 2008.

## 16.10 Exercises

1. With two independent inputs in a classification problem, that is,  $p(x_1, x_2|C) = p(x_1|C)p(x_2|C)$ , how can we calculate  $p(x_1|x_2)$ ? Derive the formula for  $p(x_j|C_i) \sim \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$ .
2. For a head-to-head node, show that equation 16.10 implies  $P(X, Y) = P(X) \cdot P(Y)$ .
3. In figure 16.4, calculate  $P(R|W)$ ,  $P(R|W, S)$ , and  $P(R|W, \sim S)$ .
4. In equation 16.30,  $X$  is binary. How do we need to modify it if  $X$  can take one of  $K$  discrete values?
5. Show that in a directed graph where the joint distribution is written as equation 16.12,  $\sum_{\mathbf{x}} p(\mathbf{x}) = 1$ .
6. Draw the Necker cube as a graphical model defining links to indicate mutually reinforcing or inhibiting relations between different corner interpretations.
7. How can we do inference on the dynamic weather graph shown in figure 16.17?
8. Write down the graphical model for linear logistic regression for two classes in the manner of figure 16.9.
9. Propose a suitable goodness measure that can be used in learning graph structure as a state-space search. What are suitable operators?
10. Generally, in a newspaper, a reporter writes a series of articles on successive days related to the same topic as the story develops. How can we model this using a graphical model?

## 16.11 References

- Andrieu, C., N. de Freitas, A. Doucet, and M. I. Jordan. 2003. "An Introduction to MCMC for Machine Learning." *Machine Learning* 50: 5-43.
- Barnard, K., P. Duygulu, D. Forsyth, N. de Freitas, D. M. Blei, and M. I. Jordan. 2003. "Matching Words and Pictures." *Journal of Machine Learning Research* 3: 1107-1135.
- Bilmes, J., and C. Bartels. 2005. "Graphical Model Architectures for Speech Recognition." *IEEE Signal Processing Magazine* 22: 89-100.
- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning*. New York: Springer.
- Buntine, W. 1996. "A Guide to the Literature on Learning Probabilistic Networks from Data." *IEEE Transactions on Knowledge and Data Engineering* 8: 195-210.

- Cowell, R. G., A. P. Dawid, S. L. Lauritzen, and D. J. Spiegelhalter. 1999. *Probabilistic Networks and Expert Systems*. New York: Springer.
- Donkers, J., and K. Tuyls. 2008. "Belief Networks in Bioinformatics." In *Computational Intelligence in Bioinformatics*, ed. A. Kelemen, A. Abraham, and Y. Chen, 75-111. Berlin: Springer.
- Ghahramani, Z. 2001. "An Introduction to Hidden Markov Models and Bayesian Networks." *International Journal of Pattern Recognition and Artificial Intelligence* 15: 9-42.
- Jensen, F. 1996. *An Introduction to Bayesian Networks*. New York: Springer.
- Jordan, M. I., ed. 1999. *Learning in Graphical Models*. Cambridge, MA: MIT Press.
- Jordan, M. I. 2004. "Graphical Models." *Statistical Science* 19: 140-155.
- Jordan, M. I. 2009. *An Introduction to Probabilistic Graphical Models*. Forthcoming.
- Jordan, M. I., Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. 1999. "An Introduction to Variational Methods for Graphical Models." In *Learning in Graphical Models*, ed. M. I. Jordan, 105-161. Cambridge, MA: MIT Press.
- Kschischang, F. R., B. J. Frey, and H.-A. Loeliger. 2001. "Factor Graphs and the Sum-Product Algorithm." *IEEE Transactions on Information Theory* 47: 498-519.
- Lauritzen, S. L., and D. J. Spiegelhalter. 1988. "Local Computations with Probabilities on Graphical Structures and their Application to Expert Systems." *Journal of Royal Statistical Society B* 50: 157-224.
- MacKay, D. J. C. 2003. *Information Theory, Inference, and Learning Algorithms*. Cambridge, UK: Cambridge University Press.
- Neapolitan, R. E. 2004. *Learning Bayesian Networks*. Upper Saddle River, NJ: Pearson.
- Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Francisco, CA: Morgan Kaufmann.
- Pearl, J. 2000. *Causality: Models, Reasoning, and Inference*. Cambridge, UK: Cambridge University Press.
- Sudderth, E. B., A. Torralba, W. T. Freeman, and A. S. Willsky. 2008. "Describing Visual Scenes Using Transformed Objects and Parts." *International Journal of Computer Vision* 77: 291-330.
- Zweig, G. 2003. "Bayesian Network Structures and Inference Techniques for Automatic Speech Recognition." *Computer Speech and Language* 17: 173-193.

# 17

## *Combining Multiple Learners*

*We discussed many different learning algorithms in the previous chapters. Though these are generally successful, no one single algorithm is always the most accurate. Now, we are going to discuss models composed of multiple learners that complement each other so that by combining them, we attain higher accuracy.*

### 17.1 Rationale

IN ANY APPLICATION, we can use one of several learning algorithms, and with certain algorithms, there are hyperparameters that affect the final learner. For example, in a classification setting, we can use a parametric classifier or a multilayer perceptron, and, for example, with a multilayer perceptron, we should also decide on the number of hidden units. The No Free Lunch Theorem states that there is no single learning algorithm that in any domain always induces the most accurate learner. The usual approach is to try many and choose the one that performs the best on a separate validation set.

Each learning algorithm dictates a certain model that comes with a set of assumptions. This inductive bias leads to error if the assumptions do not hold for the data. Learning is an ill-posed problem and with finite data, each algorithm converges to a different solution and fails under different circumstances. The performance of a learner may be fine-tuned to get the highest possible accuracy on a validation set, but this fine-tuning is a complex task and still there are instances on which even the best learner is not accurate enough. The idea is that there may be another learner that is accurate on these. By suitably combining multiple *base-learners* then, accuracy can be improved. Recently with computation and

BASE-LEARNER

memory getting cheaper, such systems composed of multiple learners have become popular (Kuncheva 2004).

There are basically two questions here:

1. How do we generate base-learners that complement each other?
2. How do we combine the outputs of base-learners for maximum accuracy?

Our discussion in this chapter will answer these two related questions. We will see that model combination is not a trick that always increases accuracy; model combination does always increase time and space complexity of training and testing, and unless base-learners are trained carefully and their decisions combined smartly, we will only pay for this extra complexity without any significant gain in accuracy.

## 17.2 Generating Diverse Learners

DIVERSITY

Since there is no point in combining learners that always make similar decisions, the aim is to be able to find a set of *diverse* learners who differ in their decisions so that they complement each other. At the same time, there cannot be a gain in overall success unless the learners are accurate, at least in their domain of expertise. We therefore have this double task of maximizing individual accuracies and the diversity between learners. Let us now discuss the different ways to achieve this.

### Different Algorithms

We can use different learning algorithms to train different base-learners. Different algorithms make different assumptions about the data and lead to different classifiers. For example, one base-learner may be parametric and another may be nonparametric. When we decide on a single algorithm, we give emphasis to a single method and ignore all others. Combining multiple learners based on multiple algorithms, we free ourselves from taking a decision and we no longer put all our eggs in one basket.

### Different Hyperparameters

We can use the same learning algorithm but use it with different hyperparameters. Examples are the number of hidden units in a multilayer

perceptron,  $k$  in  $k$ -nearest neighbor, error threshold in decision trees, the kernel function in support vector machines, and so forth. With a Gaussian parametric classifier, whether the covariance matrices are shared or not is a hyperparameter. If the optimization algorithm uses an iterative procedure such as gradient descent whose final state depends on the initial state, such as in backpropagation with multilayer perceptrons, the initial state, for example, the initial weights, is another hyperparameter. When we train multiple base-learners with different hyperparameter values, we average over this factor and reduce variance, and therefore error.

### Different Input Representations

Separate base-learners may be using different *representations* of the same input object or event, making it possible to integrate different types of sensors/measurements/modalities. Different representations make different characteristics explicit allowing better identification. In many applications, there are multiple sources of information, and it is desirable to use all of these data to extract more information and achieve higher accuracy in prediction.

#### SENSOR FUSION

For example, in speech recognition, to recognize the uttered words, in addition to the acoustic input, we can also use the video image of the speaker's lips as the words are spoken. This is similar to *sensor fusion* where the data from different sensors are integrated to extract more information for a specific application.

The simplest approach is to concatenate all data vectors and treat it as one large vector from a single source, but this does not seem theoretically appropriate since this corresponds to modeling data as sampled from one multivariate statistical distribution. Moreover, larger input dimensionalities make the systems more complex and require larger samples for the estimators to be accurate. The approach we take is to make separate predictions based on different sources using separate base-learners, then combine their predictions.

#### RANDOM SUBSPACE

Even if there is a single input representation, by choosing random subsets from it, we can have classifiers using different input features; this is called the *random subspace method* (Ho 1998). This has the effect that different learners will look at the same problem from different points of view and will be robust; it will also help reduce the curse of dimensionality because inputs are fewer dimensional.

### Different Training Sets

Another possibility is to train different base-learners by different subsets of the training set. This can be done randomly by drawing random training sets from the given sample; this is called *bagging*. Or, the learners can be trained serially so that instances on which the preceding base-learners are not accurate are given more emphasis in training later base-learners; examples are *boosting* and *cascading*, which actively try to generate complementary learners, instead of leaving this to chance.

The partitioning of the training sample can also be done based on locality in the input space so that each base-learner is trained on instances in a certain local part of the input space; this is what is done by the *mixture of experts* that we discussed in chapter 12 but that we revisit in this context of combining multiple learners. Similarly, it is possible to define the main task in terms of a number of subtasks to be implemented by the base-learners, as is done by *error-correcting output codes*.

### Diversity vs. Accuracy

One important note is that when we generate multiple base-learners, we want them to be reasonably accurate but do not require them to be very accurate individually, so they are not, and need not be, optimized separately for best accuracy. The base-learners are not chosen for their accuracy, but for their simplicity. We do require, however, that the base-learners be diverse, that is, accurate on different instances, specializing in subdomains of the problem. What we care for is the final accuracy when the base-learners are combined, rather than the accuracies of the base-learners we started from. Let us say we have a classifier that is 80 percent accurate. When we decide on a second classifier, we do not care for the overall accuracy; we care only about how accurate it is on the 20 percent that the first classifier misclassifies, as long as we know when to use which one.

This implies that the required accuracy and diversity of the learners also depend on how their decisions are to be combined, as we will discuss next. If, as in a voting scheme, a learner is consulted for all inputs, it should be accurate everywhere and diversity should be enforced everywhere; if we have a partitioning of the input space into regions of expertise for different learners, diversity is already guaranteed by this partitioning and learners need to be accurate only in their own local domains.

### 17.3 Model Combination Schemes

There are also different ways the multiple base-learners are combined to generate the final output:

#### MULTIEXPERT COMBINATION

- *Multiexpert combination* methods have base-learners that work in *parallel*. These methods can in turn be divided into two:
  - In the *global* approach, also called *learner fusion*, given an input, all base-learners generate an output and all these outputs are used. Examples are *voting* and *stacking*.
  - In the *local* approach, or *learner selection*, for example, in *mixture of experts*, there is a *gating* model, which looks at the input and chooses one (or very few) of the learners as responsible for generating the output.

#### MULTISTAGE COMBINATION

- *Multistage combination* methods use a *serial* approach where the next base-learner is trained with or tested on only the instances where the previous base-learners are not accurate enough. The idea is that the base-learners (or the different representations they use) are sorted in increasing complexity so that a complex base-learner is not used (or its complex representation is not extracted) unless the preceding simpler base-learners are not confident. An example is *cascading*.

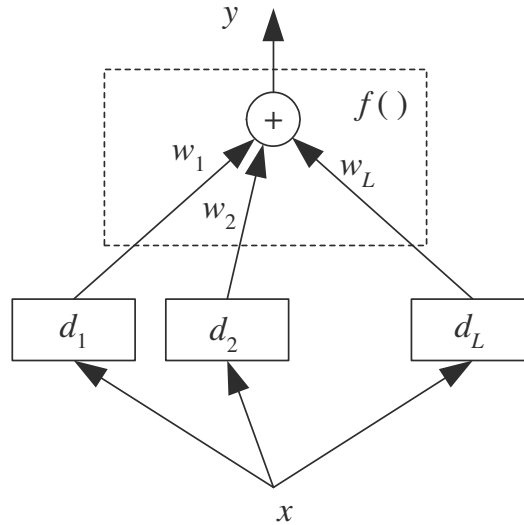
Let us say that we have  $L$  base-learners. We denote by  $d_j(x)$  the prediction of base-learner  $\mathcal{M}_j$  given the arbitrary dimensional input  $x$ . In the case of multiple representations, each  $\mathcal{M}_j$  uses a different input representation  $x_j$ . The final prediction is calculated from the predictions of the base-learners:

$$(17.1) \quad y = f(d_1, d_2, \dots, d_L | \Phi)$$

where  $f(\cdot)$  is the combining function with  $\Phi$  denoting its parameters.

When there are  $K$  outputs, for each learner there are  $d_{ji}(x), i = 1, \dots, K, j = 1, \dots, L$ , and combining them, we also generate  $K$  values,  $y_i, i = 1, \dots, K$  and then for example in classification, we choose the class with the maximum  $y_i$  value:

$$\text{Choose } C_i \text{ if } y_i = \max_{k=1}^K y_k$$



**Figure 17.1** Base-learners are  $d_j$  and their outputs are combined using  $f(\cdot)$ . This is for a single output; in the case of classification, each base-learner has  $K$  outputs that are separately used to calculate  $y_i$ , and then we choose the maximum. Note that here, all learners observe the same input; it may be the case that different learners observe different representations of the same input object or event.

## 17.4 Voting

**VOTING** The simplest way to combine multiple classifiers is by *voting*, which corresponds to taking a linear combination of the learners (see figure 17.1):

$$(17.2) \quad y_i = \sum_j w_j d_{ji} \text{ where } w_j \geq 0, \sum_j w_j = 1$$

ENSEMBLES  
LINEAR OPINION  
POOLS

This is also known as *ensembles* and *linear opinion pools*. In the simplest case, all learners are given equal weight and we have *simple voting* that corresponds to taking an average. Still, taking a (weighted) sum is only one of the possibilities and there are also other combination rules, as shown in table 17.1 (Kittler et al. 1998). If the outputs are not posterior probabilities, these rules require that outputs be normalized to the same scale (Jain, Nandakumar, and Ross 2005).

**Table 17.1** Classifier combination rules.

Rule	Fusion function $f(\cdot)$
Sum	$y_i = \frac{1}{L} \sum_{j=1}^L d_{ji}$
Weighted sum	$y_i = \sum_j w_j d_{ji}, w_j \geq 0, \sum_j w_j = 1$
Median	$y_i = \text{median}_j d_{ji}$
Minimum	$y_i = \min_j d_{ji}$
Maximum	$y_i = \max_j d_{ji}$
Product	$y_i = \prod_j d_{ji}$

**Table 17.2** Example of combination rules on three learners and three classes.

	$C_1$	$C_2$	$C_3$
$d_1$	0.2	0.5	0.3
$d_2$	0.0	0.6	0.4
$d_3$	0.4	0.4	0.2
Sum	0.2	<b>0.5</b>	0.3
Median	0.2	<b>0.5</b>	0.4
Minimum	0.0	<b>0.4</b>	0.2
Maximum	0.4	<b>0.6</b>	0.4
Product	0.0	<b>0.12</b>	0.032

An example of the use of these rules is shown in table 17.2, which demonstrates the effects of different rules. Sum rule is the most intuitive and is the most widely used in practice. Median rule is more robust to outliers; minimum and maximum rules are pessimistic and optimistic, respectively. With the product rule, each learner has veto power; regardless of the other ones, if one learner has an output of 0, the overall output goes to 0. Note that after the combination rules,  $y_i$  do not necessarily sum up to 1.

In weighted sum,  $d_{ji}$  is the vote of learner  $j$  for class  $C_i$  and  $w_j$  is the weight of its vote. Simple voting is a special case where all voters have equal weight, namely,  $w_j = 1/L$ . In classification, this is called *plurality voting* where the class having the maximum number of votes is the winner. When there are two classes, this is *majority voting* where the winning

class gets more than half of the votes (exercise 1). If the voters can also supply the additional information of how much they vote for each class (e.g., by the posterior probability), then after normalization, these can be used as weights in a *weighted voting* scheme. Equivalently, if  $d_{ji}$  are the class posterior probabilities,  $P(C_i|x, \mathcal{M}_j)$ , then we can just sum them up ( $w_j = 1/L$ ) and choose the class with maximum  $y_i$ .

In the case of regression, simple or weighted averaging or median can be used to fuse the outputs of base-regressors. Median is more robust to noise than the average.

Another possibility to find  $w_j$  is to assess the accuracies of the learners (regressor or classifier) on a separate validation set and use that information to compute the weights, so that we give more weights to more accurate learners. These weights can also be learned from data, as we will discuss when we discuss stacked generalization in section 17.9.

BAYESIAN MODEL  
COMBINATION

Voting schemes can be seen as approximations under a Bayesian framework with weights approximating prior model probabilities, and model decisions approximating model-conditional likelihoods. This is *Bayesian model combination*. For example, in classification we have  $w_j \equiv P(\mathcal{M}_j)$ ,  $d_{ji} = P(C_i|x, \mathcal{M}_j)$ , and equation 17.2 corresponds to

$$(17.3) \quad P(C_i|x) = \sum_{\text{all models } \mathcal{M}_j} P(C_i|x, \mathcal{M}_j)P(\mathcal{M}_j)$$

Simple voting corresponds to a uniform prior. If we have a prior distribution preferring simpler models, this would give larger weights to them. We cannot integrate over all models; we only choose a subset for which we believe  $P(\mathcal{M}_j)$  is high, or we can have another Bayesian step and calculate  $P(\mathcal{M}_j|X)$ , the probability of a model given the sample, and sample high probable models from this density.

Hansen and Salamon (1990) have shown that given independent two-class classifiers with success probability higher than 1/2, namely, better than random guessing, by taking a majority vote, the accuracy increases as the number of voting classifiers increases.

Let us assume that  $d_j$  are iid with expected value  $E[d_j]$  and variance  $\text{Var}(d_j)$ , then when we take a simple average with  $w_j = 1/L$ , the expected value and variance of the output are

$$E[y] = E \left[ \sum_j \frac{1}{L} d_j \right] = \frac{1}{L} L E[d_j] = E[d_j]$$

$$(17.4) \quad \text{Var}(y) = \text{Var}\left(\sum_j \frac{1}{L}d_j\right) = \frac{1}{L^2}\text{Var}\left(\sum_j d_j\right) = \frac{1}{L^2}L\text{Var}(d_j) = \frac{1}{L}\text{Var}(d_j)$$

We see that the expected value does not change, so the bias does not change. But variance, and therefore mean square error, decreases as the number of independent voters,  $L$ , increases. In the general case,

$$(17.5) \quad \text{Var}(y) = \frac{1}{L^2}\text{Var}\left(\sum_j d_j\right) = \frac{1}{L^2}\left[\sum_j \text{Var}(d_j) + 2\sum_j \sum_{i<j} \text{Cov}(d_j, d_i)\right]$$

which implies that if learners are positively correlated, variance (and error) increase. We can thus view using different algorithms and input features as efforts to decrease, if not completely eliminate, the positive correlation. In section 17.10, we will discuss pruning methods to remove learners with high positive correlation from an ensemble.

We also see here that further decrease in variance is possible if the voters are not independent but negatively correlated. The error then decreases if the accompanying increase in bias is not higher because these aims are contradictory; we cannot have a number of classifiers that are all accurate *and* negatively correlated. In mixture of experts for example, where learners are localized, the experts are negatively correlated but biased (Jacobs 1997).

If we view each base-learner as a random noise function added to the true discriminant/regression function and if these noise functions are uncorrelated with 0 mean, then the averaging of the individual estimates is like averaging over the noise. In this sense, voting has the effect of smoothing in the functional space and can be thought of as a regularizer with a smoothness assumption on the true function (Perrone 1993). We saw an example of this in figure 4.5d, where, averaging over models with large variance, we get a better fit than those of the individual models. This is the idea in voting: we vote over models with high variance and low bias so that after combination, the bias remains small and we reduce the variance by averaging. Even if the individual models are biased, the decrease in variance may offset this bias and still a decrease in error is possible.

## 17.5 Error-Correcting Output Codes

ERROR-CORRECTING  
OUTPUT CODES

In *error-correcting output codes* (ECOC) (Dietterich and Bakiri 1995), the

main classification task is defined in terms of a number of subtasks that are implemented by the base-learners. The idea is that the original task of separating one class from all other classes may be a difficult problem. Instead, we want to define a set of simpler classification problems, each specializing in one aspect of the task, and combining these simpler classifiers, we get the final classifier.

Base-learners are binary classifiers having output  $-1/+1$ , and there is a *code matrix*  $\mathbf{W}$  of  $K \times L$  whose  $K$  rows are the binary codes of classes in terms of the  $L$  base-learners  $d_j$ . For example, if the second row of  $\mathbf{W}$  is  $[-1, +1, +1, -1]$ , this means that for us to say an instance belongs to  $C_2$ , the instance should be on the negative side of  $d_1$  and  $d_4$ , and on the positive side of  $d_2$  and  $d_3$ . Similarly, the columns of the code matrix defines the task of the base-learners. For example, if the third column is  $[-1, +1, +1]^T$ , we understand that the task of the third base-learner,  $d_3$ , is to separate the instances of  $C_1$  from the instances of  $C_2$  and  $C_3$  combined. This is how we form the training set of the base-learners. For example in this case, all instances labeled with  $C_2$  and  $C_3$  form  $\mathcal{X}_3^+$  and instances labeled with  $C_1$  form  $\mathcal{X}_3^-$ , and  $d_3$  is trained so that  $x^t \in \mathcal{X}_3^+$  give output  $+1$  and  $x^t \in \mathcal{X}_3^-$  give output  $-1$ .

The code matrix thus allows us to define a polychotomy ( $K > 2$  classification problem) in terms of dichotomies ( $K = 2$  classification problem), and it is a method that is applicable using any learning algorithm to implement the dichotomizer base-learners—for example, linear or multi-layer perceptrons (with a single output), decision trees, or SVMs whose original definition is for two-class problems.

The typical one discriminant per class setting corresponds to the diagonal code matrix where  $L = K$ . For example, for  $K = 4$ , we have

$$\mathbf{W} = \begin{bmatrix} +1 & -1 & -1 & -1 \\ -1 & +1 & -1 & -1 \\ -1 & -1 & +1 & -1 \\ -1 & -1 & -1 & +1 \end{bmatrix}$$

The problem here is that if there is an error with one of the base-learners, there may be a misclassification because the class code words are so similar. So the approach in error-correcting codes is to have  $L > K$  and increase the Hamming distance between the code words. One possibility is *pairwise separation* of classes where there is a separate base-learner to separate  $C_i$  from  $C_j$ , for  $i < j$  (section 10.4). In this case,

$L = K(K - 1)/2$  and with  $K = 4$ , the code matrix is

$$\mathbf{W} = \begin{bmatrix} +1 & +1 & +1 & 0 & 0 & 0 \\ -1 & 0 & 0 & +1 & +1 & 0 \\ 0 & -1 & 0 & -1 & 0 & +1 \\ 0 & 0 & -1 & 0 & -1 & -1 \end{bmatrix}$$

where a 0 entry denotes “don’t care.” That is,  $d_1$  is trained to separate  $C_1$  from  $C_2$  and does not use the training instances belonging to the other classes. Similarly, we say that an instance belongs to  $C_2$  if  $d_1 = -1$  and  $d_4 = d_5 = +1$ , and we do not consider the values of  $d_2, d_3$ , and  $d_6$ . The problem here is that  $L$  is  $\mathcal{O}(K^2)$ , and for large  $K$  pairwise separation may not be feasible.

The approach is to set  $L$  beforehand and then find  $\mathbf{W}$  such that the distances between rows, and at the same time the distances between columns, are as large as possible, in terms of Hamming distance. With  $K$  classes, there are  $2^{(K-1)} - 1$  possible columns, namely, two-class problems. This is because  $K$  bits can be written in  $2^K$  different ways and complements (e.g., “0101” and “1010,” from our point of view, define the same discriminant) dividing the possible combinations by 2 and then subtracting 1 because a column of all 0s (or 1s) is useless. For example, when  $K = 4$ , we have

$$\mathbf{W} = \begin{bmatrix} -1 & -1 & -1 & -1 & -1 & -1 & -1 \\ -1 & -1 & -1 & +1 & +1 & +1 & +1 \\ -1 & +1 & +1 & -1 & -1 & +1 & +1 \\ +1 & -1 & +1 & -1 & +1 & -1 & +1 \end{bmatrix}$$

When  $K$  is large, for a given value of  $L$ , we look for  $L$  columns out of the  $2^{(K-1)} - 1$ . We would like these columns of  $\mathbf{W}$  to be as different as possible so that the tasks to be learned by the base-learners are as different from each other as possible. At the same time, we would like the rows of  $\mathbf{W}$  to be as different as possible so that we can have maximum error correction in case one or more base-learners fail.

ECOC can be written as a voting scheme where the entries of  $\mathbf{W}$ ,  $w_{ij}$ , are considered as vote weights:

$$(17.6) \quad y_i = \sum_{j=1}^L w_{ij} d_j$$

and then we choose the class with the highest  $y_i$ . Taking a weighted sum and then choosing the maximum instead of checking for an exact match

allows  $d_j$  to no longer need to be binary but to take a value between  $-1$  and  $+1$ , carrying soft certainties instead of hard decisions. Note that a value  $p_j$  between 0 and 1, for example, a posterior probability, can be converted to a value  $d_j$  between  $-1$  and  $+1$  simply as

$$d_j = 2p_j - 1$$

The difference between equation 17.6 and the generic voting model of equation 17.2 is that the weights of votes can be different for different classes, namely, we no longer have  $w_j$  but  $w_{ij}$ , and also that  $w_j \geq 0$  whereas  $w_{ij}$  are  $-1$ ,  $0$ , or  $+1$ .

One problem with ECOC is that because the code matrix  $\mathbf{W}$  is set a priori, there is no guarantee that the subtasks as defined by the columns of  $\mathbf{W}$  will be simple. Dietterich and Bakiri (1995) report that the dichotomizer trees may be larger than the polychotomizer trees and when multilayer perceptrons are used, there may be slower convergence by backpropagation.

## 17.6 Bagging

**BAGGING** *Bagging* is a voting method whereby base-learners are made different by training them over slightly different training sets. Generating  $L$  slightly different samples from a given sample is done by bootstrap, where given a training set  $\mathcal{X}$  of size  $N$ , we draw  $N$  instances randomly from  $\mathcal{X}$  *with replacement*. Because sampling is done with replacement, it is possible that some instances are drawn more than once and that certain instances are not drawn at all. When this is done to generate  $L$  samples  $\mathcal{X}_j, j = 1, \dots, L$ , these samples are similar because they are all drawn from the same original sample, but they are also slightly different due to chance. The base-learners  $d_j$  are trained with these  $L$  samples  $\mathcal{X}_j$ .

**UNSTABLE ALGORITHM**

A learning algorithm is an *unstable algorithm* if small changes in the training set causes a large difference in the generated learner, namely, the learning algorithm has high variance. Bagging, short for bootstrap aggregating, uses bootstrap to generate  $L$  training sets, trains  $L$  base-learners using an unstable learning procedure, and then, during testing, takes an average (Breiman 1996). Bagging can be used both for classification and regression. In the case of regression, to be more robust, one can take the median instead of the average when combining predictions.

We saw before that averaging reduces variance only if the positive correlation is small; an algorithm is stable if different runs of the same al-

gorithm on resampled versions of the same dataset lead to learners with high positive correlation. Algorithms such as decision trees and multi-layer perceptrons are unstable. Nearest neighbor is stable, but condensed nearest neighbor is unstable (Alpaydm 1997). If the original training set is large, then we may want to generate smaller sets of size  $N' < N$  from them using bootstrap, since otherwise the bootstrap replicates  $\mathcal{X}_j$  will be too similar, and  $d_j$  will be highly correlated.

## 17.7 Boosting

BOOSTING  
WEAK LEARNER  
STRONG LEARNER

In bagging, generating complementary base-learners is left to chance and to the instability of the learning method. In boosting, we actively try to generate complementary base-learners by training the next learner on the mistakes of the previous learners. The original *boosting* algorithm (Schapire 1990) combines three weak learners to generate a strong learner. A *weak learner* has error probability less than  $1/2$ , which makes it better than random guessing on a two-class problem, and a *strong learner* has arbitrarily small error probability.

Given a large training set, we randomly divide it into three. We use  $\mathcal{X}_1$  and train  $d_1$ . We then take  $\mathcal{X}_2$  and feed it to  $d_1$ . We take all instances misclassified by  $d_1$  and also as many instances on which  $d_1$  is correct from  $\mathcal{X}_2$ , and these together form the training set of  $d_2$ . We then take  $\mathcal{X}_3$  and feed it to  $d_1$  and  $d_2$ . The instances on which  $d_1$  and  $d_2$  disagree form the training set of  $d_3$ . During testing, given an instance, we give it to  $d_1$  and  $d_2$ ; if they agree, that is the response, otherwise the response of  $d_3$  is taken as the output. Schapire (1990) has shown that this overall system has reduced error rate, and the error rate can arbitrarily be reduced by using such systems recursively, that is, a boosting system of three models used as  $d_j$  in a higher system.

Though it is quite successful, the disadvantage of the original boosting method is that it requires a very large training sample. The sample should be divided into three and furthermore, the second and third classifiers are only trained on a subset on which the previous ones err. So unless one has a quite large training set,  $d_2$  and  $d_3$  will not have training sets of reasonable size. Drucker et al. (1994) use a set of 118,000 instances in boosting multilayer perceptrons for optical handwritten digit recognition.

ADABOOST Freund and Schapire (1996) proposed a variant, named *AdaBoost*, short

<p><b>Training:</b>  For all <math>\{x^t, r^t\}_{t=1}^N \in \mathcal{X}</math>, initialize <math>p_1^t = 1/N</math>  For all base-learners <math>j = 1, \dots, L</math>  Randomly draw <math>\mathcal{X}_j</math> from <math>\mathcal{X}</math> with probabilities <math>p_j^t</math>  Train <math>d_j</math> using <math>\mathcal{X}_j</math>  For each <math>(x^t, r^t)</math>, calculate <math>y_j^t \leftarrow d_j(x^t)</math>  Calculate error rate: <math>\epsilon_j \leftarrow \sum_t p_j^t \cdot 1(y_j^t \neq r^t)</math>  If <math>\epsilon_j &gt; 1/2</math>, then <math>L \leftarrow j - 1</math>; stop  <math>\beta_j \leftarrow \epsilon_j / (1 - \epsilon_j)</math>  For each <math>(x^t, r^t)</math>, decrease probabilities if correct:  If <math>y_j^t = r^t</math>, then <math>p_{j+1}^t \leftarrow \beta_j p_j^t</math> Else <math>p_{j+1}^t \leftarrow p_j^t</math>  Normalize probabilities:  <math>Z_j \leftarrow \sum_t p_{j+1}^t</math>; <math>p_{j+1}^t \leftarrow p_{j+1}^t / Z_j</math></p> <p><b>Testing:</b>  Given <math>x</math>, calculate <math>d_j(x), j = 1, \dots, L</math>  Calculate class outputs, <math>i = 1, \dots, K</math>:  <math display="block">y_i = \sum_{j=1}^L \left( \log \frac{1}{\beta_j} \right) d_{ji}(x)</math></p>
--

Figure 17.2 AdaBoost algorithm.

for adaptive boosting, that uses the same training set over and over and thus need not be large, but the classifiers should be simple so that they do not overfit. AdaBoost can also combine an arbitrary number of base-learners, not three.

Many variants of AdaBoost have been proposed; here, we discuss the original algorithm AdaBoost.M1 (see figure 17.2). The idea is to modify the probabilities of drawing the instances as a function of the error. Let us say  $p_j^t$  denotes the probability that the instance pair  $(x^t, r^t)$  is drawn to train the  $j$ th base-learner. Initially, all  $p_1^t = 1/N$ . Then we add new base-learners as follows, starting from  $j = 1$ :  $\epsilon_j$  denotes the error rate of  $d_j$ . AdaBoost requires that learners are weak, that is,  $\epsilon_j < 1/2, \forall j$ ; if not, we stop adding new base-learners. Note that this error rate is not on the original problem but on the dataset used at step  $j$ . We define  $\beta_j = \epsilon_j / (1 - \epsilon_j) < 1$ , and we set  $p_{j+1}^t = \beta_j p_j^t$  if  $d_j$  correctly classifies  $x^t$ ; otherwise,  $p_{j+1}^t = p_j^t$ . Because  $p_{j+1}^t$  should be probabilities, there is a normalization where we divide  $p_{j+1}^t$  by  $\sum_t p_{j+1}^t$ , so that they sum up to 1. This has the effect that the probability of a correctly classified instance

is decreased, and the probability of a misclassified instance increases. Then a new sample of the same size is drawn from the original sample according to these modified probabilities,  $p_{j+1}^t$ , with replacement, and is used to train  $d_{j+1}$ .

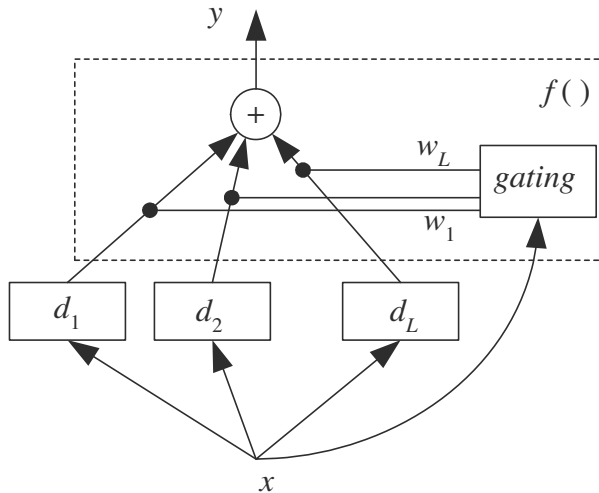
This has the effect that  $d_{j+1}$  focuses more on instances misclassified by  $d_j$ ; that is why the base-learners are chosen to be simple and not accurate, since otherwise the next training sample would contain only a few outlier and noisy instances repeated many times over. For example, with decision trees, *decision stumps*, which are trees grown only one or two levels, are used. So it is clear that these would have bias but the decrease in variance is larger and the overall error decreases. An algorithm like the linear discriminant has low variance, and we cannot gain by AdaBoosting linear discriminants.

Once training is done, AdaBoost is a voting method. Given an instance, all  $d_j$  decide and a weighted vote is taken where weights are proportional to the base-learners' accuracies (on the training set):  $w_j = \log(1/\beta_j)$ . Freund and Schapire (1996) showed improved accuracy in twenty-two benchmark problems, equal accuracy in one problem, and worse accuracy in four problems.

MARGIN Schapire et al. (1998) explain that the success of AdaBoost is due to its property of increasing the *margin*. If the margin increases, the training instances are better separated and an error is less likely. This makes AdaBoost's aim similar to that of support vector machines (chapter 13).

In AdaBoost, although different base-learners have slightly different training sets, this difference is not left to chance as in bagging, but is a function of the error of the previous base-learner. The actual performance of boosting on a particular problem is clearly dependent on the data and the base-learner. There should be enough training data and the base-learner should be weak but not too weak, and boosting is especially susceptible to noise and outliers.

AdaBoost has also been generalized to regression: One straightforward way, proposed by Avnimelech and Intrator (1997), checks for whether the prediction error is larger than a certain threshold, and if so marks it as error, then uses AdaBoost proper. In another version (Drucker 1997), probabilities are modified based on the magnitude of error, such that instances where the previous base-learner commits a large error, have a higher probability of being drawn to train the next base-learner. Weighted average, or median, is used to combine the predictions of the base-learners.



**Figure 17.3** Mixture of experts is a voting method where the votes, as given by the gating system, are a function of the input. The combiner system  $f$  also includes this gating system.

## 17.8 Mixture of Experts Revisited

### MIXTURE OF EXPERTS

In voting, the weights  $w_j$  are constant over the input space. In the *mixture of experts* architecture, which we previously discussed in section 12.8) as a local method, as an extension of radial basis functions, there is a gating network whose outputs are weights of the experts. This architecture can then be viewed as a voting method where the votes depend on the input, and may be different for different inputs. The competitive learning algorithm used by the mixture of experts localizes the base-learners such that each of them becomes an expert in a different part of the input space and have its weight,  $w_j(x)$ , close to 1 in its region of expertise. The final output is a weighted average as in voting

$$(17.7) \quad y = \sum_{j=1}^L w_j(x) d_j$$

except in this case, both the base-learners and the weights are a function of the input (see figure 17.3).

Jacobs (1997) has shown that in the mixture of experts architecture, experts are biased but are negatively correlated. As training proceeds, bias decreases and expert variances increase but at the same time as experts localize in different parts of the input space, their covariances get more and more negative, which, due to equation 17.5, decreases the total variance, and thus the error. In section 12.8, we considered the case where both experts and gating are linear functions but a nonlinear method, for example, a multilayer perceptron with hidden units, can also be used for both. This may decrease the expert biases but risks increasing expert variances and overfitting.

DYNAMIC CLASSIFIER  
SELECTION

In *dynamic classifier selection*, similar to the gating network of mixture of experts, there is first a system which takes a test input and estimates the competence of base-classifiers in the vicinity of the input. It then picks the most competent to generate output and that output is given as the overall output. Woods, Kegelmeyer, and Bowyer (1997) find the  $k$  nearest training points of the test input, look at the accuracies of the base classifiers on those, and choose the one that performs the best on them. Only the selected base-classifier need be evaluated for that test input. To decrease variance, at the expense of more computation, one can take a vote over a few competent base-classifiers instead of using just a single one.

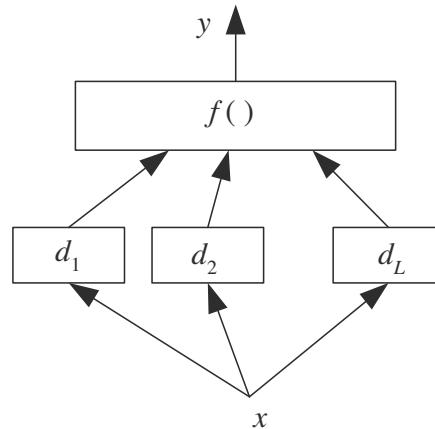
Note that in such a scheme, one should make sure that for any region of the input space, there is a competent base-classifier; this implies that there should be some partitioning of the learning of the input space among the base-classifiers. This is the nice property of mixture of experts, namely, the gating model that does the selection and the expert base-learners that it selects from are trained in a coupled manner. It would be straightforward to have a regression version of this dynamic learner selection algorithm (exercise 5).

## 17.9 Stacked Generalization

STACKED  
GENERALIZATION

*Stacked generalization* is a technique proposed by Wolpert (1992) that extends voting in that the way the output of the base-learners is combined need not be linear but is learned through a combiner system,  $f(\cdot|\Phi)$ , which is another learner, whose parameters  $\Phi$  are also trained (see figure 17.4):

$$(17.8) \quad y = f(d_1, d_2, \dots, d_L | \Phi)$$



**Figure 17.4** In stacked generalization, the combiner is another learner and is not restricted to being a linear combination as in voting.

The combiner learns what the correct output is when the base-learners give a certain output combination. We cannot train the combiner function on the training data because the base-learners may be memorizing the training set; the combiner system should actually learn how the base-learners make errors. Stacking is a means of estimating and correcting for the biases of the base-learners. Therefore, the combiner should be trained on data unused in training the base-learners.

If  $f(\cdot|w_1, \dots, w_L)$  is a linear model with constraints,  $w_i \geq 0$ ,  $\sum_j w_j = 1$ , the optimal weights can be found by constrained regression, but of course we do not need to enforce this; in stacking, there is no restriction on the combiner function and unlike voting,  $f(\cdot)$  can be nonlinear. For example, it may be implemented as a multilayer perceptron with  $\Phi$  its connection weights. The outputs of the base-learners  $d_j$  define a new  $L$ -dimensional space in which the output discriminant/regression function is learned by the combiner function.

In stacked generalization, we would like the base-learners to be as different as possible so that they will complement each other, and, for this, it is best if they are based on different learning algorithms. If we are combining classifiers that can generate continuous outputs, for example, posterior probabilities, it is better that they be the combined rather than

hard decisions.

When we compare a trained combiner as we have in stacking, with a fixed rule such as in voting, we see that both have their advantages: a trained rule is more flexible and may have less bias, but adds extra parameters, risks introducing variance, and needs extra time and data for training. Note also that there is no need to normalize classifier outputs before stacking.

## 17.10 Fine-Tuning an Ensemble

Model combination is not a magical formula always guaranteed to decrease error; base-learners should be diverse and accurate—that is, they should provide useful information. If a base-learner does not add to accuracy, it can be discarded; also, of the two base-learners that are highly correlated, one is not needed. Note that an inaccurate learner can also worsen accuracy, for example, majority voting assumes more than half of the classifiers to be accurate for an input. Therefore, given a set of candidate base-learners, it may not be a good idea to use all and we may do better by choosing a subset. This means that selecting a subset is good not only for decreasing complexity but can also improve accuracy.

### ENSEMBLE SELECTION

Choosing a subset from an ensemble of base-learners is similar to input feature selection, and the possible approaches for *ensemble selection* are the same. We can have a forward/incremental/growing approach where at each iteration, from a set of candidate base-learners, we add to the ensemble the one that most improves accuracy, we can have a backward/decremental/pruning approach where at each iteration, we remove the base-learner from the ensemble whose absence leads to highest improvement, or we can have a floating approach where both additions and removals are allowed. The combination scheme can be a fixed rule, such as voting, or it can be a trained stacker. Such a selection scheme would not include inaccurate learners, ones that are not diverse enough or are correlated (Caruana et al. 2004; Ruta and Gabrys 2005). Different learners may be using different representations, and such an approach also allows choosing the best complementary representations (Demir and Alpaydın 2005).

Actually, just as in stacking, if we consider the combination as a learner that takes base-learner outputs as inputs, what we are aiming here is input dimensionality reduction, which we discussed in chapter 6. Again,

one possibility is feature selection where we discard the uninformative inputs and keep the useful ones; in ensemble methods, this corresponds to choosing a subset from an ensemble of base-learners, as we discussed earlier. Note that if we use a decision tree as the combiner it acts both as a selector and a combiner (Ulaş et al. 2009).

The second possibility is feature extraction where from the space of the outputs of base-learners, the aim is to go to a new, lower-dimensional space where we remove unnecessary inputs and also remove correlations. Merz (1999) proposes the SCANN algorithm that uses correspondence analysis—a variant of principal components analysis (section 6.3)—on the crisp outputs of base classifiers and combines them using the nearest mean classifier. Actually, any linear or nonlinear feature extraction method we discussed in chapter 6 can be used and its (preferably continuous) output can be fed to any learner. So with  $L$  learners and  $K$  outputs each, we map from the  $K \cdot L$ -dimensional space to the new space of lower dimensional, uncorrelated space of these “eigenlearners” where we train the combiner (using a separate dataset unused to train the base-learners and the dimensionality reducer).

Rather than drastically discarding or keeping a subset of the ensemble, this approach uses all the base-learners, and hence all the information, but does not decrease complexity.

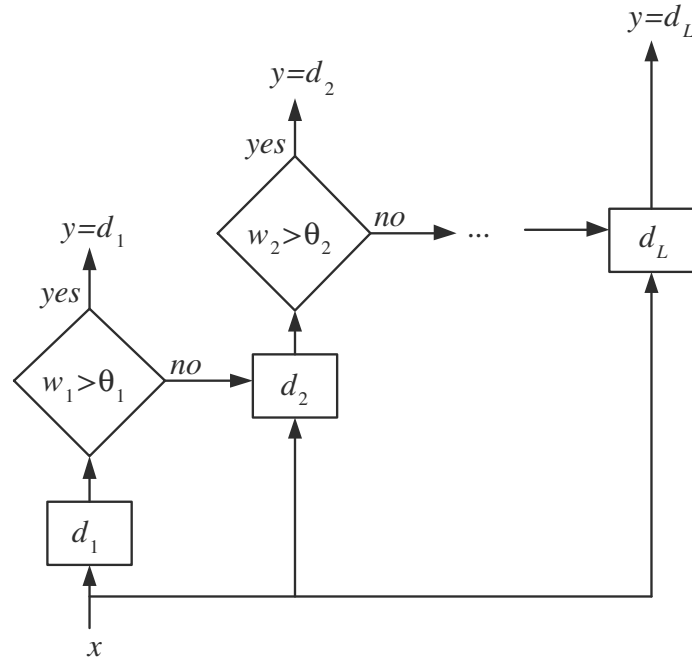
### 17.11 Cascading

CASCADING The idea in cascaded classifiers is to have a *sequence* of base-classifiers  $d_j$  sorted in terms of their space or time complexity, or the cost of the representation they use, so that  $d_{j+1}$  is costlier than  $d_j$  (Kaynak and Alpaydmn 2000). *Cascading* is a multistage method, and we use  $d_j$  only if all preceding learners,  $d_k, k < j$  are not confident (see figure 17.5). For this, associated with each learner is a *confidence*  $w_j$  such that we say  $d_j$  is confident of its output and can be used if  $w_j > \theta_j$  where  $1/K < \theta_j \leq \theta_{j+1} < 1$  is the confidence threshold. In classification, the confidence function is set to the highest posterior:  $w_j \equiv \max_i d_{ji}$ ; this is the strategy used for rejections (section 3.3).

We use learner  $d_j$  if all the preceding learners are not confident:

$$(17.9) \quad y_i = d_{ji} \text{ if } w_j > \theta_j \text{ and } \forall k < j, w_k < \theta_k$$

Starting with  $j = 1$ , given a training set, we train  $d_j$ . Then we find all instances from a separate validation set on which  $d_j$  is not confident, and



**Figure 17.5** Cascading is a multistage method where there is a sequence of classifiers, and the next one is used only when the preceding ones are not confident.

these constitute the training set of  $d_{j+1}$ . Note that unlike in AdaBoost, we choose not only the misclassified instances but the ones for which the previous base-learner is not confident. This covers the misclassifications as well as the instances for which the posterior is not high enough; these are instances on the right side of the boundary but for which the distance to the discriminant, namely, the margin, is not large enough.

The idea is that an early simple classifier handles the majority of instances, and a more complex classifier is used only for a small percentage, thereby not significantly increasing the overall complexity. This is contrary to the multiexpert methods like voting where all base-learners generate their output for any instance. If the problem space is complex, a few base-classifiers may be cascaded increasing the complexity at each stage. In order not to increase the number of base-classifiers, the few instances not covered by any are stored as they are and are treated by a

nonparametric classifier, such as  $k$ -NN.

The inductive bias of cascading is that the classes can be explained by a small number of “rules” in increasing complexity, with an additional small set of “exceptions” not covered by the rules. The rules are implemented by simple base-classifiers, for example, perceptrons of increasing complexity, which learn general rules valid over the whole input space. Exceptions are localized instances and are best handled by a nonparametric model.

Cascading thus stands between the two extremes of parametric and nonparametric classification. The former—for example, a linear model—finds a single rule that should cover all the instances. A nonparametric classifier—for example,  $k$ -NN—stores the whole set of instances without generating any simple rule explaining them. Cascading generates a rule (or rules) to explain a large part of the instances as cheaply as possible and stores the rest as exceptions. This makes sense in a lot of learning applications. For example, most of the time the past tense of a verb in English is found by adding a “-d” or “-ed” to the verb; there are also irregular verbs—for example, “go”/“went”—that do not obey this rule.

## 17.12 Notes

The idea in combining learners is to divide a complex task into simpler tasks that are handled by separately trained base-learners. Each base-learner has its own task. If we had a large learner containing all the base-learners, then it would risk overfitting. For example, consider taking a vote over three multilayer perceptrons, each with a single hidden layer. If we combine them all together with the linear model combining their outputs, this is a large multilayer perceptron with two hidden layers. If we train this large model with the whole sample, it very probably overfits. When we train the three multilayer perceptrons separately, for example, using ECOC, bagging, and so forth, it is as if we define a required output for the second-layer hidden nodes of the large multilayer perceptron. This puts a constraint on what the overall learner should learn and simplifies learning.

One disadvantage of combining is that the combined system is not interpretable. For example, even though decision trees are interpretable, bagged or boosted trees are not interpretable. Error-correcting codes with their weights as  $-1/0/+1$  allow some form of interpretability. Mayoraz

and Moreira (1997) discuss incremental methods for learning the error-correcting output codes where base-learners are added when needed. Allwein, Schapire, and Singer (2000) discuss various methods for coding multiclass problems as two-class problems. Alpaydm and Mayoraz (1999) consider the application of ECOC where linear base-learners are combined to get nonlinear discriminants, and they also propose methods to learn the ECOC matrix from data.

The earliest and most intuitive approach is voting. Kittler et al. (1998) give a review of fixed rules and also discuss an application where multiple representations are combined. The task is person identification using three representations: frontal face image, face profile image, and voice. The error rate of the voting model is lower than the error rates when a single representation is used. Another application is given in Alimođlu and Alpaydm 1997 where for improved handwritten digit recognition, two sources of information are combined: one is the temporal pen movement data as the digit is written on a touch-sensitive pad, and the other is the static two-dimensional bitmap image once the digit is written. In that application, the two classifiers using either of the two representations have around 5 percent error, but combining the two reduces the error rate to 3 percent. It is also seen that the critical stage is the design of the complementary learners and/or representations, the way they are combined is not as critical.

## BIOMETRICS

Combining different modalities is used in *biometrics*, where the aim is authentication using different input sources, fingerprint, signature, face, and so on. In such a case, different classifiers use these modalities separately and their decisions are combined. This both improves accuracy and makes *spoofing* more difficult.

Noble (2004) makes a distinction between three type of combination strategies when we have information coming from multiple sources in different representations or modalities:

- In *early integration*, all these inputs are concatenated to form a single vector that is then fed to a single classifier. Previously we discussed why this is not a very good idea.
- In *late integration*, which we advocated in this chapter, different inputs are fed to separate classifiers whose outputs are then combined, by voting, stacking, or any other method we discussed.
- Kernel algorithms, which we discussed in chapter 13, allow a different

method of integration that Noble (2004) calls *intermediate integration*, as being between early and late integration. This is the *multiple kernel learning* approach (see section 13.8) where there is a single kernel machine classifier that uses multiple kernels for different inputs and the combination is not in the input space as in early integration, or in the space of decisions as in late integration, but in the space of the basis functions that define the kernels. For different sources, there are different notions of similarity calculated by their kernels, and the classifier accumulates and uses them.

It has been shown by Jacobs (1995) that  $L$  dependent experts are worth the same as  $L'$  independent experts where  $L' \leq L$ . Under certain circumstances, voting models and Bayesian techniques will yield identical results (Jacobs 1995). The priors of equation 17.3 are in turn modeled as distributions with hyperparameters and in the ideal case, one should integrate over the whole model-parameter space. This approach is not generally feasible in practice and one resorts to approximation or sampling. With advances in Bayesian statistics, these supra-Bayesian techniques may become more important in the near future.

Combining multiple learners has been a popular topic in machine learning since the early 1990s, and research has been going on ever since. Kuncheva (2004) discusses different aspects of classifier combination; the book also includes a section on combination of multiple clustering results.

AdaBoosted decision trees used to be considered to be one of the best machine learning algorithms. There are also versions of AdaBoost where the next base-learner is trained on the residual of the previous base-learner (Hastie, Tibshirani, and Friedman 2001). Recently, it has been noticed that ensembles do not always improve accuracy and research has started to focus on the criteria that a good ensemble should satisfy or how to form a good one. A survey of the role of diversity in ensembles is given in Kuncheva 2005.

### 17.13 Exercises

1. If each base-learner is iid and correct with probability  $p > 1/2$ , what is the probability that a majority vote over  $L$  classifiers gives the correct answer?
2. In bagging, to generate the  $L$  training sets, what would be the effect of using  $L$ -fold cross-validation instead of bootstrap?

3. Propose an incremental algorithm for learning error-correcting output codes where new two-class problems are added as they are needed to better solve the multiclass problem.
4. In mixture of experts, we can have different experts use different input representations. How can we design the gating network in such a case?
5. Propose a dynamic regressor selection algorithm.
6. What is the difference between voting and stacking using a linear perceptron as the combiner function?
7. In cascading, why do we require  $\theta_{j+1} \geq \theta_j$ ?
8. To be able to use cascading for regression, during test, a regressor should be able to say if it is confident of its output. How can we implement this?
9. How can we combine the results of multiple clustering solutions?
10. In section 17.10, we discussed that if we use a decision tree as a combiner in stacking, it works both as a selector and a combiner. What are the other advantages and disadvantages?

## 17.14 References

- Alimoğlu, F., and E. Alpaydm. 1997. "Combining Multiple Representations and Classifiers for Pen-Based Handwritten Digit Recognition." In *Fourth International Conference on Document Analysis and Recognition*, 637-640. Los Alamitos, CA: IEEE Computer Society.
- Allwein, E. L., R. E. Schapire, and Y. Singer. 2000. "Reducing Multiclass to Binary: A Unifying Approach for Margin Classifiers." *Journal of Machine Learning Research* 1: 113-141.
- Alpaydm, E. 1997. "Voting over Multiple Condensed Nearest Neighbors." *Artificial Intelligence Review* 11: 115-132.
- Alpaydm, E., and E. Mayoraz. 1999. "Learning Error-Correcting Output Codes from Data." In *Ninth International Conference on Artificial Neural Networks*, 743-748. London: IEE Press.
- Avnimelech, R., and N. Intrator. 1997. "Boosting Regression Estimators." *Neural Computation* 11: 499-520.
- Breiman, L. 1996. "Bagging Predictors." *Machine Learning* 26: 123-140.
- Caruana, R., A. Niculescu-Mizil, G. Crew, and A. Ksikes. 2004. "Ensemble Selection from Libraries of Models." In *Twenty-First International Conference on Machine Learning*, ed. C. E. Brodley, 137-144. New York: ACM.
- Demir, C., and E. Alpaydm. 2005. "Cost-Conscious Classifier Ensembles." *Pattern Recognition Letters* 26: 2206-2214.

- Dietterich, T. G., and G. Bakiri. 1995. "Solving Multiclass Learning Problems via Error-Correcting Output Codes." *Journal of Artificial Intelligence Research* 2: 263-286.
- Drucker, H. 1997. "Improving Regressors using Boosting Techniques." In *Fourteenth International Conference on Machine Learning*, ed. D. H. Fisher, 107-115. San Mateo, CA: Morgan Kaufmann.
- Drucker, H., C. Cortes, L. D. Jackel, Y. Le Cun, and V. Vapnik. 1994. "Boosting and Other Ensemble Methods." *Neural Computation* 6: 1289-1301.
- Freund, Y., and R. E. Schapire. 1996. "Experiments with a New Boosting Algorithm." In *Thirteenth International Conference on Machine Learning*, ed. L. Saitta, 148-156. San Mateo, CA: Morgan Kaufmann.
- Hansen, L. K., and P. Salamon. 1990. "Neural Network Ensembles." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12: 993-1001.
- Hastie, T., R. Tibshirani, and J. Friedman. 2001. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York: Springer.
- Ho, T. K. 1998. "The Random Subspace Method for Constructing Decision Forests." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20: 832-844.
- Jacobs, R. A. 1995. "Methods for Combining Experts' Probability Assessments." *Neural Computation* 7: 867-888.
- Jacobs, R. A. 1997. "Bias/Variance Analyses for Mixtures-of-Experts Architectures." *Neural Computation* 9: 369-383.
- Jain, A., K. Nandakumar, and A. Ross. 2005. "Score Normalization in Multimodal Biometric Systems." *Pattern Recognition* 38: 2270-2285.
- Kaynak, C., and E. Alpaydm. 2000. "MultiStage Cascading of Multiple Classifiers: One Man's Noise is Another Man's Data." In *Seventeenth International Conference on Machine Learning*, ed. P. Langley, 455-462. San Francisco: Morgan Kaufmann.
- Kittler, J., M. Hatef, R. P. W. Duin, and J. Matas. 1998. "On Combining Classifiers." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20: 226-239.
- Kuncheva, L. I. 2004. *Combining Pattern Classifiers: Methods and Algorithms*. Hoboken, NJ: Wiley.
- Kuncheva, L. I. 2005. Special issue on Diversity in Multiple Classifier Systems. *Information Fusion* 6: 1-115.
- Mayoraz, E., and M. Moreira. 1997. "On the Decomposition of Polychotomies into Dichotomies." In *Fourteenth International Conference on Machine Learning*, ed. D. H. Fisher, 219-226. San Mateo, CA: Morgan Kaufmann.

- Merz, C. J. 1999. "Using Correspondence Analysis to Combine Classifiers." *Machine Learning* 36: 33-58.
- Noble, W. S. 2004. "Support Vector Machine Applications in Computational Biology." In *Kernel Methods in Computational Biology*, ed. B. Schölkopf, K. Tsuda, and J.-P. Vert, 71-92. Cambridge, MA: MIT Press.
- Perrone, M. P. 1993. "Improving Regression Estimation: Averaging Methods for Variance Reduction with Extensions to General Convex Measure." Ph.D. thesis, Brown University.
- Ruta, D., and B. Gabrys. 2005. "Classifier Selection for Majority Voting." *Information Fusion* 6: 63-81.
- Schapire, R. E. 1990. "The Strength of Weak Learnability." *Machine Learning* 5: 197-227.
- Schapire, R. E., Y. Freund, P. Bartlett, and W. S. Lee. 1998. "Boosting the Margin: A New Explanation for the Effectiveness of Voting Methods." *Annals of Statistics* 26: 1651-1686.
- Ulaş, A., M. Semerci, O. T. Yıldız, and E. Alpaydın. 2009. "Incremental Construction of Classifier and Discriminant Ensembles." *Information Sciences* 179: 1298-1318.
- Wolpert, D. H. 1992. "Stacked Generalization." *Neural Networks* 5: 241-259.
- Woods, K., W. P. Kegelmeyer Jr., and K. Bowyer. 1997. "Combination of Multiple Classifiers Using Local Accuracy Estimates." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19: 405-410.



# 18 *Reinforcement Learning*

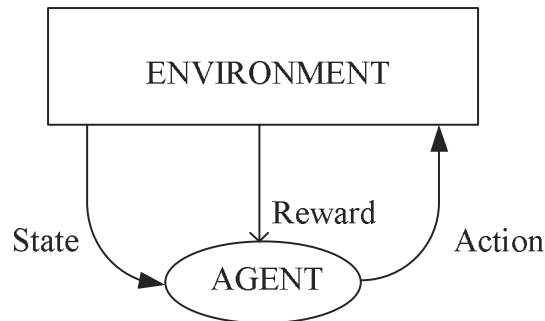
*In reinforcement learning, the learner is a decision-making agent that takes actions in an environment and receives reward (or penalty) for its actions in trying to solve a problem. After a set of trial-and-error runs, it should learn the best policy, which is the sequence of actions that maximize the total reward.*

## 18.1 Introduction

LET US SAY we want to build a machine that learns to play chess. In this case we cannot use a supervised learner for two reasons. First, it is very costly to have a teacher that will take us through many games and indicate us the best move for each position. Second, in many cases, there is no such thing as the best move; the goodness of a move depends on the moves that follow. A single move does not count; a sequence of moves is good if after playing them we win the game. The only feedback is at the end of the game when we win or lose the game.

Another example is a robot that is placed in a maze. The robot can move in one of the four compass directions and should make a sequence of movements to reach the exit. As long as the robot is in the maze, there is no feedback and the robot tries many moves until it reaches the exit and only then does it get a reward. In this case there is no opponent, but we can have a preference for shorter trajectories, implying that in this case we play against time.

These two applications have a number of points in common: there is a decision maker, called the *agent*, that is placed in an *environment* (see figure 18.1). In chess, the game-player is the decision maker and the environment is the board; in the second case, the maze is the environment



**Figure 18.1** The agent interacts with an environment. At any state of the environment, the agent takes an action that changes the state and returns a reward.

of the robot. At any time, the environment is in a certain *state* that is one of a set of possible states—for example, the state of the board, the position of the robot in the maze. The decision maker has a set of *actions* possible: legal movement of pieces on the chess board, movement of the robot in possible directions without hitting the walls, and so forth. Once an action is chosen and taken, the state changes. The solution to the task requires a sequence of actions, and we get feedback, in the form of a *reward* rarely, generally only when the complete sequence is carried out. The reward defines the problem and is necessary if we want a *learning* agent. The learning agent learns the best sequence of actions to solve a problem where “best” is quantified as the sequence of actions that has the maximum cumulative reward. Such is the setting of *reinforcement learning*.

CRITIC

CREDIT ASSIGNMENT

Reinforcement learning is different from the learning methods we discussed before in a number of respects. It is called “learning with a critic,” as opposed to learning with a teacher which we have in supervised learning. A *critic* differs from a teacher in that it does not tell us what to do but only how well we have been doing in the past; the critic never informs in advance. The feedback from the critic is scarce and when it comes, it comes late. This leads to the *credit assignment* problem. After taking several actions and getting the reward, we would like to assess the individual actions we did in the past and find the moves that led us to win the reward so that we can record and recall them later on. As we see shortly, what a reinforcement learning program does is that it learns to generate an *internal value* for the intermediate states or actions in terms of how

good they are in leading us to the goal and getting us to the real reward. Once such an internal reward mechanism is learned, the agent can just take the local actions to maximize it.

The solution to the task requires a *sequence* of actions, and from this perspective, we remember the Markov models we discussed in chapter 15. Indeed, we use a Markov decision process to model the agent. The difference is that in the case of Markov models, there is an external process that generates a sequence of signals, for example, speech, which we observe and model. In the current case, however, it is the agent that generates the sequence of actions. Previously, we also made a distinction between observable and hidden Markov models where the states are observed or hidden (and should be inferred) respectively. Similarly here, sometimes we have a partially observable Markov decision process in cases where the agent does not know its state exactly but should infer it with some uncertainty through observations using sensors. For example, in the case of a robot moving in a room, the robot may not know its exact position in the room, nor the exact location of obstacles nor the goal, and should make decisions through a limited image provided by a camera.

## 18.2 Single State Case: *K*-Armed Bandit

*K*-ARMED BANDIT

We start with a simple example. The *K-armed bandit* is a hypothetical slot machine with  $K$  levers. The action is to choose and pull one of the levers, and we win a certain amount of money that is the reward associated with the lever (action). The task is to decide which lever to pull to maximize the reward. This is a classification problem where we choose one of  $K$ . If this were supervised learning, then the teacher would tell us the correct class, namely, the lever leading to maximum earning. In this case of reinforcement learning, we can only try different levers and keep track of the best. This is a simplified reinforcement learning problem because there is only one state, or one slot machine, and we need only decide on the action. Another reason why this is simplified is that we immediately get a reward after a single action; the reward is not delayed, so we immediately see the value of our action.

Let us say  $Q(a)$  is the value of action  $a$ . Initially,  $Q(a) = 0$  for all  $a$ . When we try action  $a$ , we get reward  $r_a \geq 0$ . If rewards are deterministic, we always get the same  $r_a$  for any pull of  $a$  and in such a case, we can just set  $Q(a) = r_a$ . If we want to exploit, once we find an action  $a$  such

that  $Q(a) > 0$ , we can keep choosing it and get  $r_a$  at each pull. However, it is quite possible that there is another lever with a higher reward, so we need to explore.

We can choose different actions and store  $Q(a)$  for all  $a$ . Whenever we want to exploit, we can choose the action with the maximum value, that is,

$$(18.1) \quad \text{choose } a^* \text{ if } Q(a^*) = \max_a Q(a)$$

If rewards are not deterministic but stochastic, we get a different reward each time we choose the same action. The amount of the reward is defined by the probability distribution  $p(r|a)$ . In such a case, we define  $Q_t(a)$  as the estimate of the value of action  $a$  at time  $t$ . It is an average of all rewards received when action  $a$  was chosen before time  $t$ . An online update can be defined as

$$(18.2) \quad Q_{t+1}(a) \leftarrow Q_t(a) + \eta[r_{t+1}(a) - Q_t(a)]$$

where  $r_{t+1}(a)$  is the reward received after taking action  $a$  at time  $(t + 1)$ st time.

Note that equation 18.2 is the *delta rule* that we have used on many occasions in the previous chapters:  $\eta$  is the learning factor (gradually decreased in time for convergence),  $r_{t+1}$  is the desired output, and  $Q_t(a)$  is the current prediction.  $Q_{t+1}(a)$  is the *expected* value of action  $a$  at time  $t + 1$  and converges to the mean of  $p(r|a)$  as  $t$  increases.

The full reinforcement learning problem generalizes this simple case in a number of ways. First, we have several states. This corresponds to having several slot machines with different reward probabilities,  $p(r|s_i, a_j)$ , and we need to learn  $Q(s_i, a_j)$ , which is the value of taking action  $a_j$  when in state  $s_i$ . Second, the actions affect not only the reward but also the next state, and we move from one state to another. Third, the rewards are delayed and we need to be able to estimate immediate values from delayed rewards.

### 18.3 Elements of Reinforcement Learning

The learning decision maker is called the *agent*. The agent interacts with the *environment* that includes everything outside the agent. The agent has sensors to decide on its *state* in the environment and takes an *action* that modifies its state. When the agent takes an action, the environment

MARKOV DECISION  
PROCESS

provides a *reward*. Time is discrete as  $t = 0, 1, 2, \dots$ , and  $s_t \in S$  denotes the state of the agent at time  $t$  where  $S$  is the set of all possible states.  $a_t \in \mathcal{A}(s_t)$  denotes the action that the agent takes at time  $t$  where  $\mathcal{A}(s_t)$  is the set of possible actions in state  $s_t$ . When the agent in state  $s_t$  takes the action  $a_t$ , the clock ticks, reward  $r_{t+1} \in \mathfrak{X}$  is received, and the agent moves to the next state,  $s_{t+1}$ . The problem is modeled using a *Markov decision process* (MDP). The reward and next state are sampled from their respective probability distributions,  $p(r_{t+1}|s_t, a_t)$  and  $P(s_{t+1}|s_t, a_t)$ . Note that what we have is a *Markov* system where the state and reward in the next time step depend only on the current state and action. In some applications, reward and next state are deterministic, and for a certain state and action taken, there is one possible reward value and next state.

EPISODE  
POLICY

Depending on the application, a certain state may be designated as the initial state and in some applications, there is also an absorbing terminal (goal) state where the search ends; all actions in this terminal state transition to itself with probability 1 and without any reward. The sequence of actions from the start to the terminal state is an *episode*, or a *trial*.

The *policy*,  $\pi$ , defines the agent's behavior and is a mapping from the states of the environment to actions:  $\pi : S \rightarrow \mathcal{A}$ . The policy defines the action to be taken in any state  $s_t$ :  $a_t = \pi(s_t)$ . The *value* of a policy  $\pi$ ,  $V^\pi(s_t)$ , is the expected cumulative reward that will be received while the agent follows the policy, starting from state  $s_t$ .

FINITE-HORIZON

In the *finite-horizon* or *episodic* model, the agent tries to maximize the expected reward for the next  $T$  steps:

$$(18.3) \quad V^\pi(s_t) = E[r_{t+1} + r_{t+2} + \dots + r_{t+T}] = E \left[ \sum_{i=1}^T r_{t+i} \right]$$

INFINITE-HORIZON

Certain tasks are continuing, and there is no prior fixed limit to the episode. In the *infinite-horizon* model, there is no sequence limit, but future rewards are discounted:

$$(18.4) \quad V^\pi(s_t) = E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] = E \left[ \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i} \right]$$

DISCOUNT RATE

where  $0 \leq \gamma < 1$  is the *discount rate* to keep the return finite. If  $\gamma = 0$ , then only the immediate reward counts. As  $\gamma$  approaches 1, rewards further in the future count more, and we say that the agent becomes more farsighted.  $\gamma$  is less than 1 because there generally is a time limit to the sequence of actions needed to solve the task. The agent may be a

robot that runs on a battery. We prefer rewards sooner rather than later because we are not certain how long we will survive.

OPTIMAL POLICY For each policy  $\pi$ , there is a  $V^\pi(s_t)$ , and we want to find the *optimal policy*  $\pi^*$  such that

$$(18.5) \quad V^*(s_t) = \max_{\pi} V^\pi(s_t), \forall s_t$$

In some applications, for example, in control, instead of working with the values of states,  $V(s_t)$ , we prefer to work with the values of state-action pairs,  $Q(s_t, a_t)$ .  $V(s_t)$  denotes how good it is for the agent to be in state  $s_t$ , whereas  $Q(s_t, a_t)$  denotes how good it is to perform action  $a_t$  when in state  $s_t$ . We define  $Q^*(s_t, a_t)$  as the value, that is, the expected cumulative reward, of action  $a_t$  taken in state  $s_t$  and then obeying the optimal policy afterward. The value of a state is equal to the value of the best possible action:

$$(18.6) \quad \begin{aligned} V^*(s_t) &= \max_{a_t} Q^*(s_t, a_t) \\ &= \max_{a_t} E \left[ \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i} \right] \\ &= \max_{a_t} E \left[ r_{t+1} + \gamma \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i+1} \right] \\ &= \max_{a_t} E [r_{t+1} + \gamma V^*(s_{t+1})] \\ &= \max_{a_t} \left( E[r_{t+1}] + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, a_t) V^*(s_{t+1}) \right) \end{aligned}$$

To each possible next state  $s_{t+1}$ , we move with probability  $P(s_{t+1}|s_t, a_t)$ , and continuing from there using the optimal policy, the expected cumulative reward is  $V^*(s_{t+1})$ . We sum over all such possible next states, and we discount it because it is one time step later. Adding our immediate expected reward, we get the total expected cumulative reward for action  $a_t$ . We then choose the best of possible actions. Equation 18.6 is known as *Bellman's equation* (Bellman 1957). Similarly, we can also write

BELLMAN'S EQUATION

$$(18.7) \quad Q^*(s_t, a_t) = E[r_{t+1}] + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, a_t) \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$$

Once we have  $Q^*(s_t, a_t)$  values, we can then define our policy  $\pi$  as taking the action  $a_t^*$ , which has the highest value among all  $Q^*(s_t, a_t)$ :

$$(18.8) \quad \pi^*(s_t) : \text{Choose } a_t^* \text{ where } Q^*(s_t, a_t^*) = \max_{a_t} Q^*(s_t, a_t)$$

```

Initialize  $V(s)$  to arbitrary values
Repeat
  For all  $s \in \mathcal{S}$ 
    For all  $a \in \mathcal{A}$ 
       $Q(s, a) \leftarrow E[r|s, a] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)V(s')$ 
     $V(s) \leftarrow \max_a Q(s, a)$ 
Until  $V(s)$  converge

```

**Figure 18.2** Value iteration algorithm for model-based learning.

This means that if we have the  $Q^*(s_t, a_t)$  values, then by using a greedy search at each *local* step we get the optimal sequence of steps that maximizes the *cumulative* reward.

## 18.4 Model-Based Learning

We start with model-based learning where we completely know the environment model parameters,  $p(r_{t+1}|s_t, a_t)$  and  $P(s_{t+1}|s_t, a_t)$ . In such a case, we do not need any exploration and can directly solve for the optimal value function and policy using dynamic programming. The optimal value function is unique and is the solution to the simultaneous equations given in equation 18.6. Once we have the optimal value function, the optimal policy is to choose the action that maximizes the value in the next state:

$$(18.9) \quad \pi^*(s_t) = \arg \max_{a_t} \left( E[r_{t+1}|s_t, a_t] + \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) V^*(s_{t+1}) \right)$$

### 18.4.1 Value Iteration

VALUE ITERATION

To find the optimal policy, we can use the optimal value function, and there is an iterative algorithm called *value iteration* that has been shown to converge to the correct  $V^*$  values. Its pseudocode is given in figure 18.2.

We say that the values converged if the maximum value difference between two iterations is less than a certain threshold  $\delta$ :

$$\max_{s \in \mathcal{S}} |V^{(l+1)}(s) - V^{(l)}(s)| < \delta$$

```

Initialize a policy  $\pi'$  arbitrarily
Repeat
   $\pi \leftarrow \pi'$ 
  Compute the values using  $\pi$  by
    solving the linear equations
      
$$V^\pi(s) = E[r|s, \pi(s)] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, \pi(s)) V^\pi(s')$$

  Improve the policy at each state
      
$$\pi'(s) \leftarrow \arg \max_a (E[r|s, a] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^\pi(s'))$$

Until  $\pi = \pi'$ 

```

**Figure 18.3** Policy iteration algorithm for model-based learning.

where  $l$  is the iteration counter. Because we care only about the actions with the maximum value, it is possible that the policy converges to the optimal one even before the values converge to their optimal values. Each iteration is  $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|)$ , but frequently there is only a small number  $k < |\mathcal{S}|$  of next possible states, so complexity decreases to  $\mathcal{O}(k|\mathcal{S}||\mathcal{A}|)$ .

### 18.4.2 Policy Iteration

In policy iteration, we store and update the policy rather than doing this indirectly over the values. The pseudocode is given in figure 18.3. The idea is to start with a policy and improve it repeatedly until there is no change. The value function can be calculated by solving for the linear equations. We then check whether we can improve the policy by taking these into account. This step is guaranteed to improve the policy, and when no improvement is possible, the policy is guaranteed to be optimal. Each iteration of this algorithm takes  $\mathcal{O}(|\mathcal{A}||\mathcal{S}|^2 + |\mathcal{S}|^3)$  time that is more than that of value iteration, but policy iteration needs fewer iterations than value iteration.

## 18.5 Temporal Difference Learning

Model is defined by the reward and next state probability distributions, and as we saw in section 18.4, when we know these, we can solve for the optimal policy using dynamic programming. However, these methods are costly, and we seldom have such perfect knowledge of the environment.

The more interesting and realistic application of reinforcement learning is when we do not have the model. This requires exploration of the environment to query the model. We first discuss how this exploration is done and later see model-free learning algorithms for deterministic and nondeterministic cases. Though we are not going to assume a full knowledge of the environment model, we will however require that it be stationary.

As we will see shortly, when we explore and get to see the value of the next state and reward, we use this information to update the value of the current state. These algorithms are called *temporal difference* algorithms because what we do is look at the difference between our current estimate of the value of a state (or a state-action pair) and the discounted value of the next state and the reward received.

TEMPORAL  
DIFFERENCE

### 18.5.1 Exploration Strategies

To explore, one possibility is to use  $\epsilon$ -greedy search where with probability  $\epsilon$ , we choose one action uniformly randomly among all possible actions, namely, explore, and with probability  $1 - \epsilon$ , we choose the best action, namely, exploit. We do not want to continue exploring indefinitely but start exploiting once we do enough exploration; for this, we start with a high  $\epsilon$  value and gradually decrease it. We need to make sure that our policy is *soft*, that is, the probability of choosing any action  $a \in \mathcal{A}$  in state  $s \in S$  is greater than 0.

We can choose probabilistically, using the softmax function to convert values to probabilities

$$(18.10) \quad P(a|s) = \frac{\exp Q(s, a)}{\sum_{b \in \mathcal{A}} \exp Q(s, b)}$$

and then sample according to these probabilities. To gradually move from exploration to exploitation, we can use a “temperature” variable  $T$  and define the probability of choosing action  $a$  as

$$(18.11) \quad P(a|s) = \frac{\exp[Q(s, a)/T]}{\sum_{b \in \mathcal{A}} \exp[Q(s, b)/T]}$$

When  $T$  is large, all probabilities are equal and we have exploration. When  $T$  is small, better actions are favored. So the strategy is to start with a large  $T$  and decrease it gradually, a procedure named *annealing*, which in this case moves from exploration to exploitation smoothly in time.

### 18.5.2 Deterministic Rewards and Actions

In model-free learning, we first discuss the simpler deterministic case, where at any state-action pair, there is a single reward and next state possible. In this case, equation 18.7 reduces to

$$(18.12) \quad Q(s_t, a_t) = r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

and we simply use this as an assignment to update  $Q(s_t, a_t)$ . When in state  $s_t$ , we choose action  $a_t$  by one of the stochastic strategies we saw earlier, which returns a reward  $r_{t+1}$  and takes us to state  $s_{t+1}$ . We then update the value of *previous* action as

$$(18.13) \quad \hat{Q}(s_t, a_t) \leftarrow r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1})$$

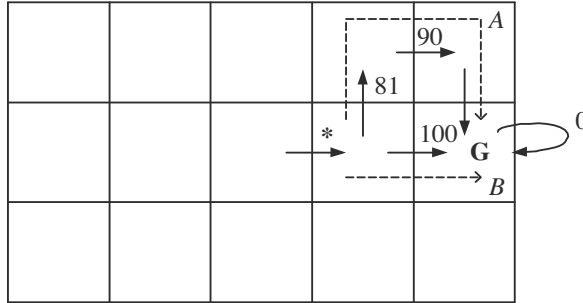
where the hat denotes that the value is an estimate.  $\hat{Q}(s_{t+1}, a_{t+1})$  is a later value and has a higher chance of being correct. We discount this by  $\gamma$  and add the immediate reward (if any) and take this as the new estimate for the previous  $\hat{Q}(s_t, a_t)$ . This is called a *backup* because it can be viewed as taking the estimated value of an action in the next time step and “backing it up” to revise the estimate for the value of a current action.

BACKUP

For now we assume that all  $\hat{Q}(s, a)$  values are stored in a table; we will see later on how we can store this information more succinctly when  $|S|$  and  $|\mathcal{A}|$  are large.

Initially all  $\hat{Q}(s_t, a_t)$  are 0, and they are updated in time as a result of trial episodes. Let us say we have a sequence of moves and at each move, we use equation 18.13 to update the estimate of the  $Q$  value of the previous state-action pair using the  $Q$  value of the current state-action pair. In the intermediate states, all rewards and therefore values are 0, so no update is done. When we get to the goal state, we get the reward  $r$  and then we can update the  $Q$  value of the previous state-action pair as  $\gamma r$ . As for the preceding state-action pair, its immediate reward is 0 and the contribution from the next state-action pair is discounted by  $\gamma$  because it is one step later. Then in another episode, if we reach this state, we can update the one preceding that as  $\gamma^2 r$ , and so on. This way, after many episodes, this information is backed up to earlier state-action pairs.  $Q$  values increase until they reach their optimal values as we find paths with higher cumulative reward, for example, shorter paths, but they never decrease (see figure 18.4).

Note that we do not know the reward or next state functions here. They are part of the environment, and it is as if we query them when



**Figure 18.4** Example to show that  $Q$  values increase but never decrease. This is a deterministic grid-world where  $G$  is the goal state with reward 100, all other immediate rewards are 0, and  $\gamma = 0.9$ . Let us consider the  $Q$  value of the transition marked by asterisk, and let us just consider only the two paths  $A$  and  $B$ . Let us say that path  $A$  is seen before path  $B$ , then we have  $\gamma \max(0, 81) = 72.9$ ; if afterward  $B$  is seen, a shorter path is found and the  $Q$  value becomes  $\gamma \max(100, 81) = 90$ . If  $B$  is seen before  $A$ , the  $Q$  value is  $\gamma \max(100, 0) = 90$ ; then when  $A$  is seen, it does not change because  $\gamma \max(100, 81) = 90$ .

we explore. We are not modeling them either, though that is another possibility. We just accept them as given and learn directly the optimal policy through the estimated value function.

**18.5.3 Nondeterministic Rewards and Actions**

If the rewards and the result of actions are not deterministic, then we have a probability distribution for the reward  $p(r_{t+1}|s_t, a_t)$  from which rewards are sampled, and there is a probability distribution for the next state  $P(s_{t+1}|s_t, a_t)$ . These help us model the uncertainty in the system that may be due to forces we cannot control in the environment: for instance, our opponent in chess, the dice in backgammon, or our lack of knowledge of the system. For example, we may have an imperfect robot which sometimes fails to go in the intended direction and deviates, or advances shorter or longer than expected.

In such a case, we have

$$(18.14) \quad Q(s_t, a_t) = E[r_{t+1}] + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, a_t) \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

We cannot do a direct assignment in this case because for the same

```

Initialize all  $Q(s, a)$  arbitrarily
For all episodes
  Initialize  $s$ 
  Repeat
    Choose  $a$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
    Take action  $a$ , observe  $r$  and  $s'$ 
    Update  $Q(s, a)$ :
       $Q(s, a) \leftarrow Q(s, a) + \eta(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$ 
     $s \leftarrow s'$ 
  Until  $s$  is terminal state

```

**Figure 18.5**  $Q$  learning, which is an off-policy temporal difference algorithm.

state and action, we may receive different rewards or move to different next states. What we do is keep a running average. This is known as the  $Q$  learning algorithm:

$Q$  LEARNING

$$(18.15) \quad \hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \eta(r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1}) - \hat{Q}(s_t, a_t))$$

We think of  $r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1})$  values as a sample of instances for each  $(s_t, a_t)$  pair and we would like  $\hat{Q}(s_t, a_t)$  to converge to its mean. As usual  $\eta$  is gradually decreased in time for convergence, and it has been shown that this algorithm converges to the optimal  $Q^*$  values (Watkins and Dayan 1992). The pseudocode of the  $Q$  learning algorithm is given in figure 18.5.

We can also think of equation 18.15 as reducing the difference between the current  $Q$  value and the backed-up estimate, from one time step later. Such algorithms are called *temporal difference* (TD) algorithms (Sutton 1988).

TEMPORAL  
DIFFERENCE

OFF-POLICY  
ON-POLICY

SARSA

This is an *off-policy* method as the value of the best next action is used without using the policy. In an *on-policy* method, the policy is used to determine also the next action. The on-policy version of  $Q$  learning is the *Sarsa* algorithm whose pseudocode is given in figure 18.6. We see that instead of looking for all possible next actions  $a'$  and choosing the best, the on-policy Sarsa uses the policy derived from  $Q$  values to choose one next action  $a'$  and uses its  $Q$  value to calculate the temporal difference. On-policy methods estimate the value of a policy while using it to take actions. In off-policy methods, these are separated, and the policy used to generate behavior, called the *behavior* policy, may in fact be differ-

```

Initialize all  $Q(s, a)$  arbitrarily
For all episodes
  Initialize  $s$ 
  Choose  $a$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
  Repeat
    Take action  $a$ , observe  $r$  and  $s'$ 
    Choose  $a'$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
    Update  $Q(s, a)$ :
       $Q(s, a) \leftarrow Q(s, a) + \eta(r + \gamma Q(s', a') - Q(s, a))$ 
       $s \leftarrow s', a \leftarrow a'$ 
  Until  $s$  is terminal state

```

**Figure 18.6** Sarsa algorithm, which is an on-policy version of  $Q$  learning.

ent from the policy that is evaluated and improved, called the *estimation* policy.

Sarsa converges with probability 1 to the optimal policy and state-action values if a *GLIE policy* is employed to choose actions. A GLIE (greedy in the limit with infinite exploration) policy is where (1) all state-action pairs are visited an infinite number of times, and (2) the policy converges in the limit to the greedy policy (which can be arranged, e.g., with  $\epsilon$ -greedy policies by setting  $\epsilon = 1/t$ ).

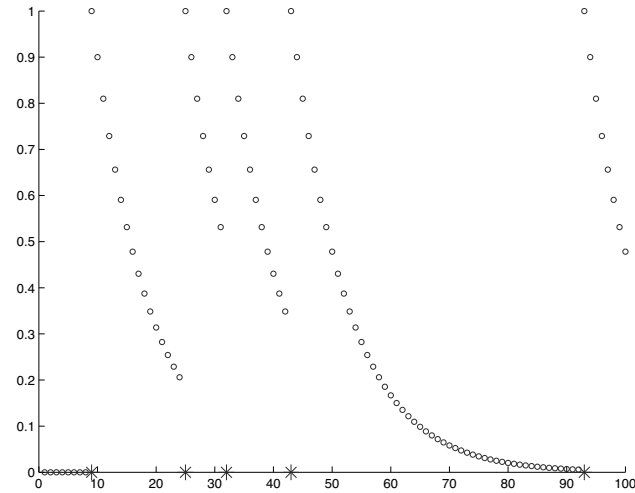
TD LEARNING The same idea of temporal difference can also be used to learn  $V(s)$  values, instead of  $Q(s, a)$ . *TD learning* (Sutton 1988) uses the following update rule to update a state value:

$$(18.16) \quad V(s_t) \leftarrow V(s_t) + \eta[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

This again is the delta rule where  $r_{t+1} + \gamma V(s_{t+1})$  is the better, later prediction and  $V(s_t)$  is the current estimate. Their difference is the temporal difference, and the update is done to decrease this difference. The update factor  $\eta$  is gradually decreased, and TD is guaranteed to converge to the optimal value function  $V^*(s)$ .

#### 18.5.4 Eligibility Traces

ELIGIBILITY TRACE The previous algorithms are one-step—that is, the temporal difference is used to update only the previous value (of the state or state-action pair). An *eligibility trace* is a record of the occurrence of past visits that en-



**Figure 18.7** Example of an eligibility trace for a value. Visits are marked by an asterisk.

ables us to implement temporal credit assignment, allowing us to update the values of previously occurring visits as well. We discuss how this is done with Sarsa to learn  $Q$  values; adapting this to learn  $V$  values is straightforward.

To store the eligibility trace, we require an additional memory variable associated with each state-action pair,  $e(s, a)$ , initialized to 0. When the state-action pair  $(s, a)$  is visited, namely, when we take action  $a$  in state  $s$ , its eligibility is set to 1; the eligibilities of all other state-action pairs are multiplied by  $\gamma\lambda$ .  $0 \leq \lambda \leq 1$  is the trace decay parameter.

$$(18.17) \quad e_t(s, a) = \begin{cases} 1 & \text{if } s = s_t \text{ and } a = a_t, \\ \gamma\lambda e_{t-1}(s, a) & \text{otherwise} \end{cases}$$

If a state-action pair has never been visited, its eligibility remains 0; if it has been, as time passes and other state-actions are visited, its eligibility decays depending on the value of  $\gamma$  and  $\lambda$  (see figure 18.7).

We remember that in Sarsa, the temporal error at time  $t$  is

$$(18.18) \quad \delta_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

In Sarsa with an eligibility trace, named Sarsa( $\lambda$ ), *all* state-action pairs

```

Initialize all  $Q(s, a)$  arbitrarily,  $e(s, a) \leftarrow 0, \forall s, a$ 
For all episodes
  Initialize  $s$ 
  Choose  $a$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
  Repeat
    Take action  $a$ , observe  $r$  and  $s'$ 
    Choose  $a'$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
     $\delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$ 
     $e(s, a) \leftarrow 1$ 
    For all  $s, a$ :
       $Q(s, a) \leftarrow Q(s, a) + \eta \delta e(s, a)$ 
       $e(s, a) \leftarrow \gamma \lambda e(s, a)$ 
     $s \leftarrow s', a \leftarrow a'$ 
  Until  $s$  is terminal state

```

**Figure 18.8** Sarsa( $\lambda$ ) algorithm.

are updated as

$$(18.19) \quad Q(s, a) \leftarrow Q(s, a) + \eta \delta_t e_t(s, a), \quad \forall s, a$$

This updates all eligible state-action pairs, where the update depends on how far they have occurred in the past. The value of  $\lambda$  defines the temporal credit: if  $\lambda = 0$ , only a one-step update is done. The algorithms we discussed in section 18.5.3 are such, and for this reason they are named  $Q(0)$ , Sarsa(0), or TD(0). As  $\lambda$  gets closer to 1, more of the previous steps are considered. When  $\lambda = 1$ , all previous steps are updated and the credit given to them falls only by  $\gamma$  per step. In online updating, all eligible values are updated immediately after each step; in offline updating, the updates are accumulated and a single update is done at the end of the episode. Online updating takes more time but converges faster. The pseudocode for *Sarsa*( $\lambda$ ) is given in figure 18.8.  $Q(\lambda)$  and TD( $\lambda$ ) algorithms can similarly be derived (Sutton and Barto 1998).

SARSA( $\lambda$ )

## 18.6 Generalization

Until now, we assumed that the  $Q(s, a)$  values (or  $V(s)$ , if we are estimating values of states) are stored in a lookup table, and the algorithms

we considered earlier are called *tabular* algorithms. There are a number of problems with this approach: (1) when the number of states and the number of actions is large, the size of the table may become quite large; (2) states and actions may be continuous, for example, turning the steering wheel by a certain angle, and to use a table, they should be discretized which may cause error; and (3) when the search space is large, too many episodes may be needed to fill in all the entries of the table with acceptable accuracy.

Instead of storing the  $Q$  values as they are, we can consider this a regression problem. This is a supervised learning problem where we define a regressor  $Q(s, a | \theta)$ , taking  $s$  and  $a$  as inputs and parameterized by a vector of parameters,  $\theta$ , to learn  $Q$  values. For example, this can be an artificial neural network with  $s$  and  $a$  as its inputs, one output, and  $\theta$  its connection weights.

A good function approximator has the usual advantages and solves the problems discussed previously. A good approximation may be achieved with a simple model without explicitly storing the training instances; it can use continuous inputs; and it allows generalization. If we know that similar  $(s, a)$  pairs have similar  $Q$  values, we can generalize from past cases and come up with good  $Q(s, a)$  values even if that state-action pair has never been encountered before.

To be able to train the regressor, we need a training set. In the case of Sarsa(0), we saw before that we would like  $Q(s_t, a_t)$  to get close to  $r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$ . So, we can form a set of training samples where the input is the state-action pair  $(s_t, a_t)$  and the required output is  $r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$ . We can write the squared error as

$$(18.20) \quad E^t(\theta) = [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]^2$$

Training sets can similarly be defined for  $Q(0)$  and TD(0), where in the latter case we learn  $V(s)$ , and the required output is  $r_{t+1} - \gamma V(s_{t+1})$ . Once such a set is ready, we can use any supervised learning algorithm for learning the training set.

If we are using a gradient-descent method, as in training neural networks, the parameter vector is updated as

$$(18.21) \quad \Delta \theta = \eta [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)] \nabla_{\theta} Q(s_t, a_t)$$

This is a one-step update. In the case of Sarsa( $\lambda$ ), the eligibility trace is also taken into account:

$$(18.22) \quad \Delta \theta = \eta \delta_t \mathbf{e}_t$$

where the temporal difference error is

$$\delta_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

and the vector of eligibilities of parameters are updated as

$$(18.23) \quad \mathbf{e}_t = \gamma \lambda \mathbf{e}_{t-1} + \nabla_{\boldsymbol{\theta}} Q(s_t, a_t)$$

with  $\mathbf{e}_0$  all zeros. In the case of a tabular algorithm, the eligibilities are stored for the state-action pairs because they are the parameters (stored as a table). In the case of an estimator, eligibility is associated with the parameters of the estimator. We also note that this is very similar to the momentum method for stabilizing backpropagation (section 11.8.1). The difference is that in the case of momentum previous weight changes are remembered, whereas here previous gradient vectors are remembered. Depending on the model used for  $Q(s_t, a_t)$ , for example, a neural network, we plug its gradient vector in equation 18.23.

In theory, any regression method can be used to train the  $Q$  function, but the particular task has a number of requirements. First, it should allow generalization; that is, we really need to guarantee that similar states and actions have similar  $Q$  values. This also requires a good coding of  $s$  and  $a$ , as in any application, to make the similarities apparent. Second, reinforcement learning updates provide instances one by one and not as a whole training set, and the learning algorithm should be able to do individual updates to learn the new instance without forgetting what has been learned before. For example, a multilayer perceptron using backpropagation can be trained with a single instance only if a small learning rate is used. Or, such instances may be collected to form a training set and learned altogether but this slows down learning as no learning happens while a sufficiently large sample is being collected.

Because of these reasons, it seems a good idea to use local learners to learn the  $Q$  values. In such methods, for example, radial basis functions, information is localized and when a new instance is learned, only a local part of the learner is updated without possibly corrupting the information in another part. The same requirements apply if we are estimating the state values as  $V(s_t | \boldsymbol{\theta})$ .

## 18.7 Partially Observable States

### 18.7.1 The Setting

In certain applications, the agent does not know the state exactly. It is equipped with sensors that return an *observation*, which the agent then uses to estimate the state. Let us say we have a robot that navigates in a room. The robot may not know its exact location in the room, or what else is there in the room. The robot may have a camera with which sensory observations are recorded. This does not tell the robot its state exactly but gives some indication as to its likely state. For example, the robot may only know that there is an obstacle to its right.

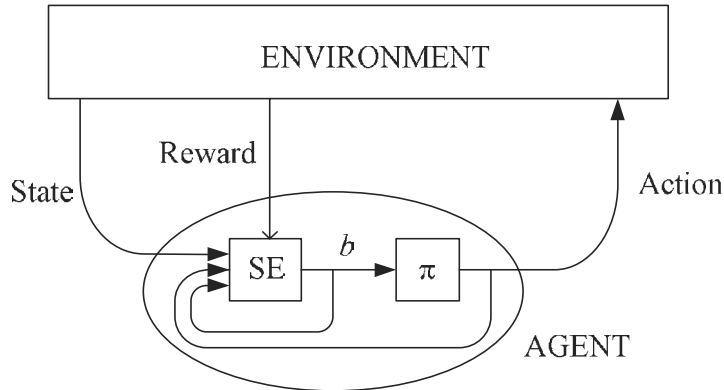
The setting is like a Markov decision process, except that after taking an action  $a_t$ , the new state  $s_{t+1}$  is not known, but we have an observation  $o_{t+1}$  that is a stochastic function of  $s_t$  and  $a_t$ :  $p(o_{t+1}|s_t, a_t)$ . This is called a *partially observable MDP* (POMDP). If  $o_{t+1} = s_{t+1}$ , then POMDP reduces to the MDP. This is just like the distinction between observable and hidden Markov models and the solution is similar; that is, from the observation, we need to infer the state (or rather a probability distribution for the states) and then act based on this. If the agent believes that it is in state  $s_1$  with probability 0.4 and in state  $s_2$  with probability 0.6, then the value of any action is 0.4 times the value of the action in  $s_1$  plus 0.6 times the value of the action in  $s_2$ .

The Markov property does not hold for observations. The next state observation does not only depend on the current action and observation. When there is limited observation, two states may appear the same but are different and if these two states require different actions, this can lead to a loss of performance, as measured by the cumulative reward. The agent should somehow compress the past trajectory into a current unique state estimate. These past observations can also be taken into account by taking a past window of observations as input to the policy, or one can use a recurrent neural network (section 11.12.2) to maintain the state without forgetting past observations.

At any time, the agent may calculate the most likely state and take an action accordingly. Or it may take an action to gather information and reduce uncertainty, for example, search for a landmark, or stop to ask for direction. This implies the importance of the *value of information*, and indeed POMDPs can be modeled as *dynamic* influence diagrams (section 16.8). The agent chooses between actions based on the amount of

PARTIALLY  
OBSERVABLE MDP

VALUE OF  
INFORMATION



**Figure 18.9** In the case of a partially observable environment, the agent has a state estimator (SE) that keeps an internal belief state  $b$  and the policy  $\pi$  generates actions based on the belief states.

information they provide, the amount of reward they produce, and how they change the state of the environment.

BELIEF STATE

To keep the process Markov, the agent keeps an internal *belief state*  $b_t$  that summarizes its experience (see figure 18.9). The agent has a *state estimator* that updates the belief state  $b_{t+1}$  based on the last action  $a_t$ , current observation  $o_{t+1}$ , and its previous belief state  $b_t$ . There is a policy  $\pi$  that generates the next action  $a_{t+1}$  based on this belief state, as opposed to the actual state that we had in a completely observable environment. The belief state is a probability distribution over states of the environment given the initial belief state (before we did any actions) and the past observation-action history of the agent (without leaving out any information that could improve agent's performance).  $Q$  learning in such a case involves the belief state-action pair values, instead of the actual state-action pairs:

$$(18.24) \quad Q(b_t, a_t) = E[r_{t+1}] + \gamma \sum_{b_{t+1}} P(b_{t+1}|b_t, a_t) V(b_{t+1})$$

### 18.7.2 Example: The Tiger Problem

We now discuss an example that is a slightly different version of the *Tiger problem* discussed in Kaelbling, Littman, and Cassandra 1998, modified as in the example in Thrun, Burgard, and Fox 2005. Let us say we are

standing in front of two doors, one to our left and the other to other right, leading to two rooms. Behind one of the two doors, we do not know which, there is a crouching tiger, and behind the other, there is a treasure. If we open the door of the room where the tiger is, we get a large negative reward, and if we open the door of the treasure room, we get some positive reward. The hidden state,  $z_L$ , is the location of the tiger. Let us say  $p$  denotes the probability that tiger is in the room to the left and therefore, the tiger is in the room to the right with probability  $1 - p$ :

$$p \equiv P(z_L = 1)$$

The two actions are  $a_L$  and  $a_R$ , which respectively correspond to opening the left or the right door. The rewards are

$r(A, Z)$	Tiger left	Tiger right
Open left	-100	+80
Open right	+90	-100

We can calculate the expected reward for the two actions. There are no future rewards because the episode ends once we open one of the doors.

$$R(a_L) = r(a_L, z_L)P(z_L) + r(a_L, z_R)P(z_R) = -100p + 80(1 - p)$$

$$R(a_R) = r(a_R, z_L)P(z_L) + r(a_R, z_R)P(z_R) = 90p - 100(1 - p)$$

Given these rewards, if  $p$  is close to 1, if we believe that there is a high chance that the tiger is on the left, the right action will be to choose the right door, and, similarly, for  $p$  close to 0, it is better to choose the left door.

The two intersect for  $p$  around 0.5, and there the expected reward is approximately  $-10$ . The fact that the expected reward is negative when  $p$  is around 0.5 (when we have uncertainty) indicates the importance of collecting information. If we can add sensors to decrease uncertainty—that is, move  $p$  away from 0.5 to either close to 0 or close to 1—we can take actions with high positive rewards. That sensing action,  $a_S$ , may have a small negative reward:  $R(a_S) = -1$ ; this may be considered as the cost of sensing or equivalent to discounting future reward by  $\gamma < 1$  because we are postponing taking the real action (of opening one of the doors).

In such a case, the expected rewards and value of the best action are shown in figure 18.10a:

$$V = \max(a_L, a_R, a_S)$$

Let us say as sensory input, we use microphones to check whether the tiger is behind the left or the right door. But we have unreliable sensors (so that we still stay in the realm of partial observability). Let us say we can only detect tiger's presence with 0.7 probability:

$$\begin{aligned} P(o_L|z_L) &= 0.7 & P(o_L|z_R) &= 0.3 \\ P(o_R|z_L) &= 0.3 & P(o_R|z_R) &= 0.7 \end{aligned}$$

If we sense  $o_L$ , our belief in the tiger's position changes:

$$p' = P(z_L|o_L) = \frac{P(o_L|z_L)P(z_L)}{p(o_L)} = \frac{0.7p}{0.7p + 0.3(1-p)}$$

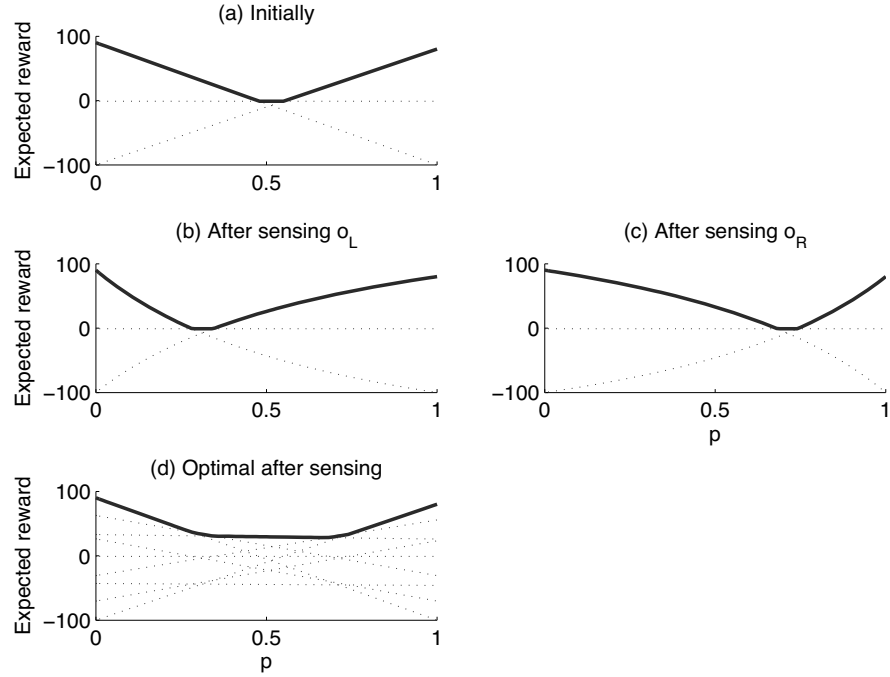
The effect of this is shown in figure 18.10b where we plot  $R(a_L|o_L)$ . Sensing  $o_L$  turns opening the right door into a better action for a wider range. The better sensors we have (if the probability of correct sensing moves from 0.7 closer to 1), the larger this range gets (exercise 9). Similarly, as we see in figure 18.10c, if we sense  $o_R$ , this increases the chances of opening the left door. Note that sensing also decreases the range where there is a need to sense (once more).

The expected rewards for the actions in this case are

$$\begin{aligned} R(a_L|o_L) &= r(a_L, z_L)P(z_L|o_L) + r(a_L, z_R)P(z_R|o_L) \\ &= -100p' + 80(1-p') \\ &= -100 \cdot \frac{0.7 \cdot p}{p(o_L)} + 80 \cdot \frac{0.3 \cdot (1-p)}{p(o_L)} \\ R(a_R|o_L) &= r(a_R, z_L)P(z_L|o_L) + r(a_R, z_R)P(z_R|o_L) \\ &= 90p' - 100(1-p') \\ &= 90 \cdot \frac{0.7 \cdot p}{p(o_L)} - 100 \cdot \frac{0.3 \cdot (1-p)}{p(o_L)} \\ R(a_S|o_L) &= -1 \end{aligned}$$

The best action in this case is the maximum of these three. Similarly, if we sense  $o_R$ , the expected rewards become

$$\begin{aligned} R(a_L|o_R) &= r(a_L, z_L)P(z_L|o_R) + r(a_L, z_R)P(z_R|o_R) \\ &= -100 \cdot \frac{0.3 \cdot p}{p(o_R)} + 80 \cdot \frac{0.7 \cdot (1-p)}{p(o_R)} \end{aligned}$$



**Figure 18.10** Expected rewards and the effect of sensing in the Tiger problem.

$$\begin{aligned}
 R(a_R|o_R) &= r(a_R, z_L)P(z_L|o_R) + r(a_R, z_R)P(z_R|o_R) \\
 &= 90 \cdot \frac{0.3 \cdot p}{p(o_R)} - 100 \cdot \frac{0.7 \cdot (1-p)}{p(o_R)} \\
 R(a_S|o_R) &= -1
 \end{aligned}$$

To calculate the expected reward, we need to take average over both sensor readings weighted by their probabilities:

$$\begin{aligned}
 V' &= \sum_j \left[ \max_i R(a_i|o_j) \right] P(o_j) \\
 &= \max(R(a_L|o_L), R(a_R|o_L), R(a_S|o_L))P(o_L) + \\
 &\quad \max(R(a_L|o_R), R(a_R|o_R), R(a_S|o_R))P(o_R) \\
 &= \max(-70p + 24(1-p), 63p - 30(1-p), -0.7p - 0.3(1-p)) + \\
 &\quad \max(-30p + 56(1-p), 27p - 70(1-p), -0.3p - 0.7(1-p))
 \end{aligned}$$

$$(18.25) \quad = \max \begin{pmatrix} -100p & +80(1-p) \\ -43p & -46(1-p) \\ 33p & +26(1-p) \\ 90p & -100(1-p) \end{pmatrix}$$

Note that when we multiply by  $P(o_L)$ , it cancels out and we get functions linear in  $p$ . These five lines and the piecewise function that corresponds to their maximum are shown in figure 18.10d. Note that the line,  $-40p - 5(1-p)$ , as well as the ones involving  $a_S$ , are beneath others for all values of  $p$  and can safely be pruned. The fact that figure 18.10d is better than figure 18.10a indicates the *value of information*.

VALUE OF  
INFORMATION

What we calculate here is the value of the best action had we chosen  $a_S$ . For example, the first line corresponds to choosing  $a_L$  after  $a_S$ . So to find the best decision with an episode of length two, we need to back this up by subtracting  $-1$ , which is the reward of  $a_S$ , and get the expected reward for the action of sense. Equivalently, we can consider this as waiting that has an immediate reward of 0 but discounts the future reward by some  $\gamma < 1$ . We also have the two usual actions of  $a_L$  and  $a_R$  and we choose the best of three; the two immediate actions and the one discounted future action.

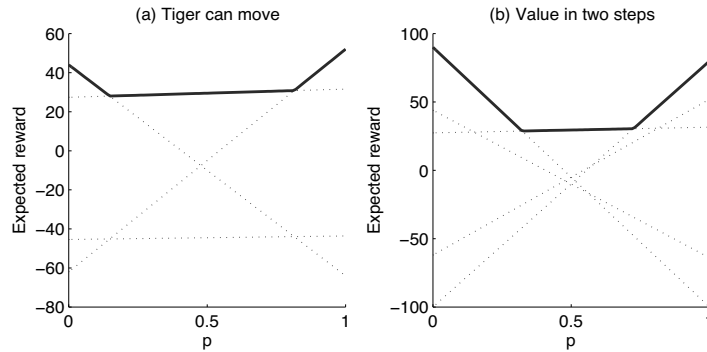
Let us now make the problem more interesting, as in the example of Thrun, Burgard, and Fox 2005. Let us assume that there is a door between the two rooms and without us seeing, the tiger can move from one room to the other. Let us say that this is a restless tiger and it stays in the same room with probability 0.2 and moves to the other room with probability 0.8. This means that  $p$  should also be updated as

$$p' = 0.2p + 0.8(1-p)$$

and this updated  $p$  should be used in equation 18.25 while choosing the best action after having chosen  $a_S$ :

$$V' = \max \begin{pmatrix} -100p' & +80(1-p') \\ 33p' & +26(1-p') \\ 90p' & -100(1-p') \end{pmatrix}$$

Figure 18.11b corresponds to figure 18.10d with the updated  $p'$ . Now, when planning for episodes of length two, we have the two immediate



**Figure 18.11** Expected rewards change (a) if the hidden state can change, and (b) when we consider episodes of length two.

actions of  $a_L$  and  $a_R$ , or we wait and sense when  $p$  changes and then we take the action and get its discounted reward (figure 18.11b):

$$V_2 = \max \begin{pmatrix} -100p & +80(1-p) \\ 90p & -100(1-p) \\ \max V' - 1 \end{pmatrix}$$

We see that figure 18.11b is better than figure 18.10a; when wrong actions may lead to large penalty, it is better to defer judgment, look for extra information, and plan ahead. We can consider longer episodes by continuing the iterative updating of  $p$  and discounting by subtracting 1 and including the two immediate actions to calculate  $V_t, t > 2$ .

The algorithm we have just discussed where the value is represented by piecewise linear functions works only when the number of states, actions, observations, and the episode length are all finite. Even in applications where any of these is not small, or when any is continuous-valued, the complexity becomes high and we need to resort to approximate algorithms having reasonable complexity. Reviews of such algorithms are given in Hauskrecht 2000 and Thrun, Burgard, and Fox 2005.

## 18.8 Notes

More information on reinforcement learning can be found in the textbook by Sutton and Barto (1998) that discusses all the aspects, learning algorithms, and several applications. A comprehensive tutorial is Kaelbling,

Littman, and Moore 1996. Recent work on reinforcement learning applied to robotics with some impressive applications is given in Thrun, Burgard, and Fox 2005.

Dynamic programming methods are discussed in Bertsekas 1987 and in Bertsekas and Tsitsiklis 1996, and TD( $\lambda$ ) and  $Q$ -learning can be seen as stochastic approximations to dynamic programming (Jaakkola, Jordan, and Singh 1994). Reinforcement learning has two advantages over classical dynamic programming: first, as they learn, they can focus on the parts of the space that are important and ignore the rest; and second, they can employ function approximation methods to represent knowledge that allows them to generalize and learn faster.

LEARNING AUTOMATA

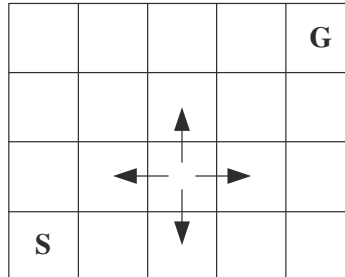
A related field is that of *learning automata* (Narendra and Thathachar 1974), which are finite state machines that learn by trial and error for solving problems like the  $K$ -armed bandit. The setting we have here is also the topic of optimal control where there is a controller (agent) taking actions in a plant (environment) that minimize cost (maximize reward).

The earliest use of temporal difference method was in Samuel's checkers program written in 1959 (Sutton and Barto 1998). For every two successive positions in a game, the two board states are evaluated by the board evaluation function that then causes an update to decrease the difference. There has been much work on games because games are both easily defined and challenging. A game like chess can easily be simulated: the allowed moves are formal, and the goal is well defined. Despite the simplicity of defining the game, expert play is quite difficult.

TD-GAMMON

One of the most impressive application of reinforcement learning is the *TD-Gammon* program that learns to play backgammon by playing against itself (Tesauro 1995). This program is superior to the previous neurogammon program also developed by Tesauro, which was trained in a supervised manner based on plays by experts. Backgammon is a complex task with approximately  $10^{20}$  states, and there is randomness due to the roll of dice. Using the TD( $\lambda$ ) algorithm, the program achieves master level play after playing 1,500,000 games against a copy of itself.

Another interesting application is in *job shop scheduling*, or finding a schedule of tasks satisfying temporal and resource constraints (Zhang and Dietterich 1996). Some tasks have to be finished before others can be started, and two tasks requiring the same resource cannot be done simultaneously. Zhang and Dietterich used reinforcement learning to quickly find schedules that satisfy the constraints and are short. Each state is one schedule, actions are schedule modifications, and the program finds not



**Figure 18.12** The grid world. The agent can move in the four compass directions starting from  $S$ . The goal state is  $G$ .

only one good schedule but a schedule for a class of related scheduling problems.

Recently hierarchical methods have also been proposed where the problem is decomposed into a set of subproblems. This has the advantage that policies learned for the subproblems can be shared for multiple problems, which accelerates learning a new problem (Dietterich 2000). Each subproblem is simpler and learning them separately is faster. The disadvantage is that when they are combined, the policy may be suboptimal.

Though reinforcement learning algorithms are slower than supervised learning algorithms, it is clear that they have a wider variety of application and have the potential to construct better learning machines (Ballard 1997). They do not need any supervision, and this may actually be better since then they are not biased by the teacher. For example, Tesauro's TD-Gammon program in certain circumstances came up with moves that turned out to be superior to those made by the best players. The field of reinforcement learning is developing rapidly, and we may expect to see other impressive results in the near future.

## 18.9 Exercises

1. Given the grid world in figure 18.12, if the reward on reaching on the goal is 100 and  $\gamma = 0.9$ , calculate manually  $Q^*(s, a)$ ,  $V^*(S)$ , and the actions of optimal policy.
2. With the same configuration given in exercise 1, use  $Q$  learning to learn the

optimal policy.

3. In exercise 1, how does the optimal policy change if another goal state is added to the lower-right corner? What happens if a state of reward  $-100$  (a very bad state) is defined in the lower-right corner?
4. Instead of having  $\gamma < 1$ , we can have  $\gamma = 1$  but with a negative reward of  $-c$  for all intermediate (nongoal) states. What is the difference?
5. In exercise 1, assume that the reward on arrival to the goal state is normal distributed with mean 100 and variance 40. Assume also that the actions are also stochastic in that when the robot advances in a direction, it moves in the intended direction with probability 0.5 and there is a 0.25 probability that it moves in one of the lateral directions. Learn  $Q(s, a)$  in this case.
6. Assume we are estimating the value function for states  $V(s)$  and that we want to use TD( $\lambda$ ) algorithm. Derive the tabular value iteration update.
7. Using equation 18.22, derive the weight update equations when a multilayer perceptron is used to estimate  $Q$ .
8. Give an example of a reinforcement learning application that can be modeled by a POMDP. Define the states, actions, observations, and reward.
9. In the tiger example, show that as we get a more reliable sensor, the range where we need to sense once again decreases.
10. Rework the tiger example using the following reward matrix

$r(A, Z)$	Tiger left	Tiger right
Open left	-100	+10
Open right	20	-100

## 18.10 References

- Ballard, D. H. 1997. *An Introduction to Natural Computation*. Cambridge, MA: MIT Press.
- Bellman, R. E. 1957. *Dynamic Programming*. Princeton: Princeton University Press.
- Bertsekas, D. P. 1987. *Dynamic Programming: Deterministic and Stochastic Models*. New York: Prentice Hall.
- Bertsekas, D. P., and J. N. Tsitsiklis. 1996. *Neuro-Dynamic Programming*. Belmont, MA: Athena Scientific.
- Dietterich, T. G. 2000. "Hierarchical Reinforcement Learning with the MAXQ Value Decomposition." *Journal of Artificial Intelligence Research* 13: 227-303.

- Hauskrecht, M. 2000. "Value-Function Approximations for Partially Observable Markov Decision Processes." *Journal of Artificial Intelligence Research* 13: 33-94.
- Jaakkola, T., M. I. Jordan, and S. P. Singh. 1994. "On the Convergence of Stochastic Iterative Dynamic Programming Algorithms." *Neural Computation* 6: 1185-1201.
- Kaelbling, L. P., M. L. Littman, and A. R. Cassandra. 1998. "Planning and Acting in Partially Observable Stochastic Domains." *Artificial Intelligence* 101: 99-134.
- Kaelbling, L. P., M. L. Littman, and A. W. Moore. 1996. "Reinforcement Learning: A Survey." *Journal of Artificial Intelligence Research* 4: 237-285.
- Narendra, K. S., and M. A. L. Thathachar. 1974. "Learning Automata—A Survey." *IEEE Transactions on Systems, Man, and Cybernetics* 4: 323-334.
- Sutton, R. S. 1988. "Learning to Predict by the Method of Temporal Differences." *Machine Learning* 3: 9-44.
- Sutton, R. S., and A. G. Barto. 1998. *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press.
- Tesauro, G. 1995. "Temporal Difference Learning and TD-Gammon." *Communications of the ACM* 38(3): 58-68.
- Thrun, S., W. Burgard, and D. Fox. 2005. *Probabilistic Robotics*. Cambridge, MA: MIT Press.
- Watkins, C. J. C. H., and P. Dayan. 1992. "Q-learning." *Machine Learning* 8: 279-292.
- Zhang, W., and T. G. Dietterich. 1996. "High-Performance Job-Shop Scheduling with a Time-Delay TD( $\lambda$ ) Network." In *Advances in Neural Information Processing Systems 8*, ed. D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, 1024-1030. Cambridge, MA: The MIT Press.

# 19 *Design and Analysis of Machine Learning Experiments*

*We discuss the design of machine learning experiments to assess and compare the performances of learning algorithms in practice and the statistical tests to analyze the results of these experiments.*

## 19.1 Introduction

IN PREVIOUS chapters, we discussed several learning algorithms and saw that, given a certain application, more than one is applicable. Now, we are concerned with two questions:

1. How can we assess the expected error of a learning algorithm on a problem? That is, for example, having used a classification algorithm to train a classifier on a dataset drawn from some application, can we say with enough confidence that later on when it is used in real life, its expected error rate will be less than, for example, 2 percent?
2. Given two learning algorithms, how can we say one has less error than the other one, for a given application? The algorithms compared can be different, for example, parametric versus nonparametric, or they can use different hyperparameter settings. For example, given a multi-layer perceptron (chapter 11) with four hidden units and another one with eight hidden units, we would like to be able to say which one has less expected error. Or with the  $k$ -nearest neighbor classifier (chapter 8), we would like to find the best value of  $k$ .

We cannot look at the training set errors and decide based on those. The error rate on the training set, by definition, is always smaller than the error rate on a test set containing instances unseen during training.

Similarly, training errors cannot be used to compare two algorithms. This is because over the training set, the more complex model having more parameters will almost always give fewer errors than the simple one.

So as we have repeatedly discussed, we need a validation set that is different from the training set. Even over a validation set though, just one run may not be enough. There are two reasons for this: First, the training and validation sets may be small and may contain exceptional instances, like noise and outliers, which may mislead us. Second, the learning method may depend on other random factors affecting generalization. For example, with a multilayer perceptron trained using backpropagation, because gradient descent converges to the nearest local minimum, the initial weights affect the final weights, and given the exact same architecture and training set, starting from different initial weights, there may be multiple possible final classifiers having different error rates on the same validation set. We thus would like to have several runs to average over such sources of randomness. If we train and validate only once, we cannot test for the effect of such factors; this is only admissible if the learning method is so costly that it can be trained and validated only once.

We use a *learning algorithm* on a dataset and generate a *learner*. If we do the training once, we have one learner and one validation error. To average over randomness (in training data, initial weights, etc.), we use the same algorithm and generate multiple learners. We test them on multiple validation sets and record a sample of validation errors. (Of course, all the training and validation sets should be drawn from the same application.) We base our evaluation of the learning algorithm on the *distribution* of these validation errors. We can use this distribution for assessing the *expected error* of the learning algorithm for that problem, or compare it with the error rate distribution of some other learning algorithm.

EXPECTED ERROR

Before proceeding to how this is done, it is important to stress a number of points:

1. We should keep in mind that whatever conclusion we draw from our analysis is conditioned on the dataset we are given. We are not comparing learning algorithms in a domain independent way but on some particular application. We are not saying anything about the expected error of a learning algorithm, or comparing one learning algorithm with another algorithm, in general. Any result we have is only true for the particular application, and only insofar as that application is rep-

NO FREE LUNCH  
THEOREM

resented in the sample we have. And anyway, as stated by the *No Free Lunch Theorem* (Wolpert 1995), there is no such thing as the “best” learning algorithm. For any learning algorithm, there is a dataset where it is very accurate and another dataset where it is very poor. When we say that a learning algorithm is good, we only quantify how well its inductive bias matches the properties of the data.

2. The division of a given dataset into a number of training and validation set pairs is only for testing purposes. Once all the tests are complete and we have made our decision as to the final method or hyperparameters, to train the final learner, we can use all the labeled data that we have previously used for training or validation.
3. Because we also use the validation set(s) for testing purposes, for example, for choosing the better of two learning algorithms, or to decide where to stop learning, it effectively becomes part of the data we use. When after all such tests, we decide on a particular algorithm and want to report its expected error, we should use a separate *test set* for this purpose, unused during training this final system. This data should have never been used before for training or validation and should be large for the error estimate to be meaningful. So, given a dataset, we should first leave some part of it aside as the test set and use the rest for training and validation. Typically, we can leave one-third of the sample as the test set, then use two-thirds for cross-validation to generate multiple training/validation set pairs, as we will see shortly. So, the training set is used to optimize the parameters, given a particular learning algorithm and model structure; the validation set is used to optimize the hyperparameters of the learning algorithm or the model structure; and the test set is used at the end, once both these have been optimized. For example, with an MLP, the training set is used to optimize the weights, the validation set is used to decide on the number of hidden units, how long to train, the learning rate, and so forth. Once the best MLP configuration is chosen, its final error is calculated on the test set. With  $k$ -NN, the training set is stored as the lookup table; we optimize the distance measure and  $k$  on the validation set and test finally on the test set.
4. In general, we compare learning algorithms by their error rates, but it should be kept in mind that in real life, error is only one of the criteria that affect our decision. Some other criteria are (Turney 2000):

- risks when errors are generalized using loss functions, instead of 0/1 loss (section 3.3),
- training time and space complexity,
- testing time and space complexity,
- interpretability, namely, whether the method allows knowledge extraction which can be checked and validated by experts, and
- easy programmability.

COST-SENSITIVE  
LEARNING

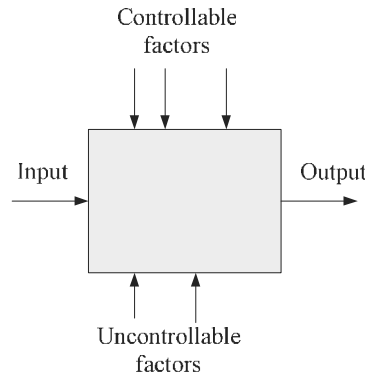
The relative importances of these factors change depending on the application. For example, if the training is to be done once in the factory, then training time and space complexity are not important; if adaptability during use is required, then they do become important. Most of the learning algorithms use 0/1 loss and take error as the single criterion to be minimized; recently, *cost-sensitive learning* variants of these algorithms have also been proposed to take other cost criteria into account.

When we train a learner on a dataset using a training set and test its accuracy on some validation set and try to draw conclusions, what we are doing is experimentation. Statistics defines a methodology to design experiments correctly and analyze the collected data in a manner so as to be able to extract significant conclusions (Montgomery 2005). In this chapter, we will see how this methodology can be used in the context of machine learning.

## 19.2 Factors, Response, and Strategy of Experimentation

EXPERIMENT

As in other branches of science and engineering, in machine learning too, we do experiments to get information about the process under scrutiny. In our case, this is a learner, which, having been trained on a dataset, generates an output for a given input. An *experiment* is a test or a series of tests where we play with the *factors* that affect the output. These factors may be the algorithm used, the training set, input features, and so on, and we observe the changes in the *response* to be able to extract information. The aim may be to identify the most important factors, screen the unimportant ones, or find the configuration of the factors that optimizes the response—for example, classification accuracy on a given test set.



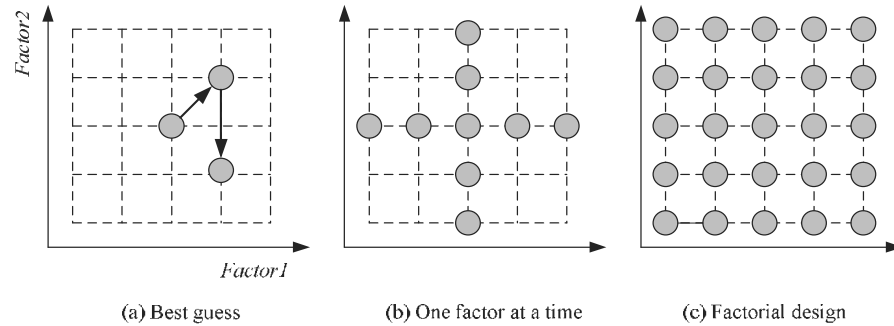
**Figure 19.1** The process generates an output given an input and is affected by controllable and uncontrollable factors.

Our aim is to plan and conduct machine learning experiments and analyze the data resulting from the experiments, to be able to eliminate the effect of chance and obtain conclusions which we can consider *statistically significant*. In machine learning, we target a learner having the highest generalization accuracy and the minimal complexity (so that its implementation is cheap in time and space) and is robust, that is, minimally affected by external sources of variability.

A trained learner can be shown as in figure 19.1; it gives an output, for example, a class code for a test input, and this depends on two type of factors: The *controllable factors*, as the name suggests, are those we have control on. The most basic is the learning algorithm used. There are also the hyperparameters of the algorithm, for example, the number of hidden units for a multilayer perceptron,  $k$  for  $k$ -nearest neighbor,  $C$  for support vector machines, and so on. The dataset used and the input representation, that is, how the input is coded as a vector, are other controllable factors.

There are also *uncontrollable factors* over which we have no control, adding undesired variability to the process, which we do not want to affect our decisions. Among these are the noise in the data, the particular training subset if we are resampling from a large set, randomness in the optimization process, for example, the initial state in gradient descent with multilayer perceptrons, and so on.

We use the output to generate the *response* variable—for example, av-



**Figure 19.2** Different strategies of experimentation with two factors and five levels each.

erage classification error on a test set, or the expected risk using a loss function, or some other measure, such as precision and recall, as we will discuss shortly.

Given several factors, we need to find the best setting for best response, or in the general case, determine their effect on the response variable. For example, we may be using principal components analyzer (PCA) to reduce dimensionality to  $d$  before a  $k$ -nearest neighbor ( $k$ -NN) classifier. The two factors are  $d$  and  $k$ , and the question is to decide which combination of  $d$  and  $k$  leads to highest performance. Or, we may be using a support vector machine classifier with Gaussian kernel, and we have the regularization parameter  $C$  and the spread of the Gaussian  $s^2$  to fine-tune together.

#### STRATEGIES OF EXPERIMENTATION

There are several *strategies of experimentation*, as shown in figure 19.2. In the *best guess* approach, we start at some setting of the factors that we believe is a good configuration. We test the response there and we fiddle with the factors one (or very few) at a time, testing each combination until we get to a state that we consider is good enough. If the experimenter has a good intuition of the process, this may work well; but note that there is no systematic approach to modify the factors and when we stop, we have no guarantee of finding the best configuration.

Another strategy is to modify *one factor at a time* where we decide on a baseline (default) value for all factors, and then we try different levels for one factor while keeping all other factors at their baseline. The major disadvantage of this is that it assumes that there is no *interaction* between the factors, which may not always be true. In the PCA/ $k$ -NN cascade we discussed earlier, each choice for  $d$  defines a different input

FACTORIAL DESIGN

space for  $k$ -NN where a different  $k$  value may be appropriate.

The correct approach is to use a *factorial design* where factors are varied together, instead of one at a time; this is colloquially called *grid search*. With  $F$  factors at  $L$  levels each, searching one factor at a time takes  $\mathcal{O}(L \cdot F)$  time, whereas a factorial experiment takes  $\mathcal{O}(L^F)$  time.

### 19.3 Response Surface Design

To decrease the number of runs necessary, one possibility is to run a fractional factorial design where we run only a subset, another is to try to use knowledge gathered from previous runs to estimate configurations that seem likely to have high response. In searching one factor at a time, if we can assume that the response is typically quadratic (with a single maximum, assuming we are maximizing a response value, such as the test accuracy), then instead of trying all values, we can have an iterative procedure where starting from some initial runs, we fit a quadratic, find its maximum analytically, take that as the next estimate, run an experiment there, add the resulting data to the sample, and then continue fitting and sampling, until we get no further improvement.

RESPONSE SURFACE DESIGN

With many factors, this is generalized as the *response surface design* method where we try to fit a parametric response function to the factors as

$$r = g(f_1, f_2, \dots, f_F | \phi)$$

where  $r$  is the response and  $f_i, i = 1, \dots, F$  are the factors. This fitted parametric function defined given the parameters  $\phi$  is our empirical model estimating the response for a particular configuration of the (controllable) factors; the effect of uncontrollable factors is modeled as noise.  $g(\cdot)$  is a (typically quadratic) regression model and after a small number of runs around some baseline (as defined by a so-called *design matrix*), one can have enough data to fit  $g(\cdot)$  on. Then, we can analytically calculate the values of  $f_i$  where the fitted  $g$  is maximum, which we take as our next guess, run an experiment there, get a data instance, add it to the sample, fit  $g$  once more, and so on, until there is convergence. Whether this approach will work well or not depends on whether the response can indeed be written as a quadratic function of the factors with a single maximum.

## 19.4 Randomization, Replication, and Blocking

Let us now talk about the three basic principles of experimental design.

- RANDOMIZATION ■ *Randomization* requires that the order in which the runs are carried out should be randomly determined so that the results are independent. This is typically a problem in real-world experiments involving physical objects; for example, machines require some time to warm up until they operate in their normal range so tests should be done in random order for time not to bias the results. Ordering generally is not a problem in software experiments.
  
- REPLICATION ■ *Replication* implies that for the same configuration of (controllable) factors, the experiment should be run a number of times to average over the effect of uncontrollable factors. In machine learning, this is typically done by running the same algorithm on a number of resampled versions of the same dataset; this is known as *cross-validation*, which we will discuss in section 19.6. How the response varies on these different replications of the same experiment allows us to obtain an estimate of the experimental error (the effect of uncontrollable factors), which we can in turn use to determine how large differences should be to be deemed *statistically significant*.
  
- BLOCKING ■ *Blocking* is used to reduce or eliminate the variability due to *nuisance factors* that influence the response but in which we are not interested. For example, defects produced in a factory may also depend on the different batches of raw material, and this effect should be isolated from the controllable factors in the factory, such as the equipment, personnel, and so on. In machine learning experimentation, when we use resampling and use different subsets of the data for different replicates, we need to make sure that for example if we are comparing learning algorithms, they should all use the same set of resampled subsets, otherwise the differences in accuracies would depend not only on the algorithms but also on the different subsets—to be able to measure the difference due to algorithms only, the different training sets in replicated runs should be identical; this is what we mean by blocking.
  
- PAIRING ■ In statistics, if there are two populations, this is called *pairing* and is used in *paired testing*.

## 19.5 Guidelines for Machine Learning Experiments

Before we start experimentation, we need to have a good idea about what it is we are studying, how the data is to be collected, and how we are planning to analyze it. The steps in machine learning are the same as for any type of experimentation (Montgomery 2005). Note that at this point, it is not important whether the task is classification or regression, or whether it is an unsupervised or a reinforcement learning application. The same overall discussion applies; the difference is only in the sampling distribution of the response data that is collected.

### A. Aim of the Study

We need to start by stating the problem clearly, defining what the objectives are. In machine learning, there may be several possibilities. As we discussed before, we may be interested in assessing the expected error (or some other response measure) of a learning algorithm on a particular problem and check that, for example, the error is lower than a certain acceptable level.

Given two learning algorithms and a particular problem as defined by a dataset, we may want to determine which one has less generalization error. These can be two different algorithms, or one can be a proposed improvement of the other, for example, by using a better feature extractor.

In the general case, we may have more than two learning algorithms, and we may want to choose the one with the least error, or order them in terms of error, for a given dataset.

In an even more general setting, instead of on a single dataset, we may want to compare two or more algorithms on two or more datasets.

### B. Selection of the Response Variable

We need to decide on what we should use as the quality measure. Most frequently, error is used that is the misclassification error for classification and mean square error for regression. We may also use some variant; for example, generalizing from 0/1 to an arbitrary loss, we may use a risk measure. In information retrieval, we use measures such as precision and recall; we will discuss such measures in section 19.7. In a cost-sensitive

setting, not only the output but also system parameters, for example, its complexity, are taken into account.

### C. Choice of Factors and Levels

What the factors are depend on the aim of the study. If we fix an algorithm and want to find the best hyperparameters, then those are the factors. If we are comparing algorithms, the learning algorithm is a factor. If we have different datasets, they also become a factor.

The levels of a factor should be carefully chosen so as not to miss a good configuration and avoid doing unnecessary experimentation. It is always good to try to normalize factor levels. For example, in optimizing  $k$  of  $k$ -nearest neighbor, one can try values such as 1, 3, 5, and so on, but in optimizing the spread  $h$  of Parzen windows, we should not try absolute values such as 1.0, 2.0, and so on, because that depends on the scale of the input; it is better to find some statistic that is an indicator of scale—for example, the average distance between an instance and its nearest neighbor—and try  $h$  as different multiples of that statistic.

Though previous expertise is a plus in general, it is also important to investigate all factors and factor levels that may be of importance and not be overly influenced by past experience.

### D. Choice of Experimental Design

It is always better to do a factorial design unless we are sure that the factors do not interact, because mostly they do. Replication number depends on the dataset size; it can be kept small when the dataset is large; we will discuss this in the next section when we talk about resampling. However, too few replicates generate few data and this will make comparing distributions difficult; in the particular case of parametric tests, the assumptions of Gaussianity may not be tenable.

Generally, given some dataset, we leave some part as the test set and use the rest for training and validation, probably many times by resampling. How this division is done is important. In practice, using small datasets leads to responses with high variance, and the differences will not be significant and results will not be conclusive.

It is also important to avoid as much as possible toy, synthetic data and use datasets that are collected from real-world under real-life circumstances. Didactic one- or two-dimensional datasets may help provide

intuition, but the behavior of the algorithms may be completely different in high-dimensional spaces.

### **E. Performing the Experiment**

Before running a large factorial experiment with many factors and levels, it is best if one does a few trial runs for some random settings to check that all is as expected. In a large experiment, it is always a good idea to save intermediate results (or seeds of the random number generator), so that a part of the whole experiment can be rerun when desired. All the results should be reproducible. In running a large experiment with many factors and factor levels, one should be aware of the possible negative effects of software aging.

It is important that an experimenter be unbiased during experimentation. In comparing one's favorite algorithm with a competitor, both should be investigated equally diligently. In large-scale studies, it may even be envisaged that testers be different from developers.

One should avoid the temptation to write one's own "library" and instead, as much as possible, use code from reliable sources; such code would have been better tested and optimized.

As in any software development study, the advantages of good documentation cannot be underestimated, especially when working in groups. All the methods developed for high-quality software engineering should also be used in machine learning experiments.

### **F. Statistical Analysis of the Data**

This corresponds to analyzing data in a way so that whatever conclusion we get is not subjective or due to chance. We cast the questions that we want to answer in a hypothesis testing framework and check whether the sample supports the hypothesis. For example, the question "Is  $A$  a more accurate algorithm than  $B$ ?" becomes the hypothesis "Can we say that the average error of learners trained by  $A$  is significantly lower than the average error of learners trained by  $B$ ?"

As always, visual analysis is helpful, and we can use histograms of error distributions, whisker-and-box plots, range plots, and so on.

## G. Conclusions and Recommendations

Once all data is collected and analyzed, we can draw objective conclusions. One frequently encountered conclusion is the need for further experimentation. Most statistical, and hence machine learning or data mining, studies are iterative. It is for this reason that we never start with all the experimentation. It is suggested that no more than 25 percent of the available resources should be invested in the first experiment (Montgomery 2005). The first runs are for investigation only. That is also why it is a good idea not to start with high expectations, or promises to one's boss or thesis advisor.

We should always remember that statistical testing never tells us if the hypothesis is correct or false, but how much the sample seems to concur with the hypothesis. There is always a risk that we do not have a conclusive result or that our conclusions be wrong, especially if the data is small and noisy.

When our expectations are not met, it is most helpful to investigate why they are not. For example, in checking why our favorite algorithm  $A$  has worked awfully bad on some cases, we can get a splendid idea for some improved version of  $A$ . All improvements are due to the deficiencies of the previous version; finding a deficiency is but a helpful hint that there is an improvement we can make!

But we should not go to the next step of testing the improved version before we are sure that we have completely analyzed the current data and learned all we could learn from it. Ideas are cheap, and useless unless tested, which is costly.

## 19.6 Cross-Validation and Resampling Methods

For replication purposes, our first need is to get a number of training and validation set pairs from a dataset  $\mathcal{X}$  (after having left out some part as the test set). To get them, if the sample  $\mathcal{X}$  is large enough, we can randomly divide it into  $K$  parts, then randomly divide each part into two and use one half for training and the other half for validation.  $K$  is typically 10 or 30. Unfortunately, datasets are never large enough to do this. So we should do our best with small datasets. This is done by repeated use of the same data split differently; this is called *cross-validation*. The catch is that this makes the error percentages dependent as these different sets share data.

CROSS-VALIDATION

STRATIFICATION

So, given a dataset  $\mathcal{X}$ , we would like to generate  $K$  training/validation set pairs,  $\{\mathcal{T}_i, \mathcal{V}_i\}_{i=1}^K$ , from this dataset. We would like to keep the training and validation sets as large as possible so that the error estimates are robust, and at the same time, we would like to keep the overlap between different sets as small as possible. We also need to make sure that classes are represented in the right proportions when subsets of data are held out, not to disturb the class prior probabilities; this is called *stratification*. If a class has 20 percent examples in the whole dataset, in all samples drawn from the dataset, it should also have approximately 20 percent examples.

### 19.6.1 K-Fold Cross-Validation

K-FOLD  
CROSS-VALIDATION

In *K-fold cross-validation*, the dataset  $\mathcal{X}$  is divided randomly into  $K$  equal-sized parts,  $\mathcal{X}_i, i = 1, \dots, K$ . To generate each pair, we keep one of the  $K$  parts out as the validation set and combine the remaining  $K - 1$  parts to form the training set. Doing this  $K$  times, each time leaving out another one of the  $K$  parts out, we get  $K$  pairs:

$$\begin{aligned} \mathcal{V}_1 &= \mathcal{X}_1 & \mathcal{T}_1 &= \mathcal{X}_2 \cup \mathcal{X}_3 \cup \dots \cup \mathcal{X}_K \\ \mathcal{V}_2 &= \mathcal{X}_2 & \mathcal{T}_2 &= \mathcal{X}_1 \cup \mathcal{X}_3 \cup \dots \cup \mathcal{X}_K \\ & \vdots & & \\ \mathcal{V}_K &= \mathcal{X}_K & \mathcal{T}_K &= \mathcal{X}_1 \cup \mathcal{X}_2 \cup \dots \cup \mathcal{X}_{K-1} \end{aligned}$$

There are two problems with this. First, to keep the training set large, we allow validation sets that are small. Second, the training sets overlap considerably, namely, any two training sets share  $K - 2$  parts.

LEAVE-ONE-OUT

$K$  is typically 10 or 30. As  $K$  increases, the percentage of training instances increases and we get more robust estimators, but the validation set becomes smaller. Furthermore, there is the cost of training the classifier  $K$  times, which increases as  $K$  is increased. As  $N$  increases,  $K$  can be smaller; if  $N$  is small,  $K$  should be large to allow large enough training sets. One extreme case of  $K$ -fold cross-validation is *leave-one-out* where given a dataset of  $N$  instances, only one instance is left out as the validation set (instance) and training uses the  $N - 1$  instances. We then get  $N$  separate pairs by leaving out a different instance at each iteration. This is typically used in applications such as medical diagnosis, where labeled data is hard to find. Leave-one-out does not permit stratification.

Recently, with computation getting cheaper, it has also become possible to have multiple runs of  $K$ -fold cross-validation, for example,  $10 \times 10$ -

fold, and use average over averages to get more reliable error estimates (Bouckaert 2003).

### 19.6.2 5×2 Cross-Validation

5 × 2  
CROSS-VALIDATION

Dietterich (1998) proposed the 5 × 2 *cross-validation*, which uses training and validation sets of equal size. We divide the dataset  $\mathcal{X}$  randomly into two parts,  $\mathcal{X}_1^{(1)}$  and  $\mathcal{X}_1^{(2)}$ , which gives our first pair of training and validation sets,  $\mathcal{T}_1 = \mathcal{X}_1^{(1)}$  and  $\mathcal{V}_1 = \mathcal{X}_1^{(2)}$ . Then we swap the role of the two halves and get the second pair:  $\mathcal{T}_2 = \mathcal{X}_1^{(2)}$  and  $\mathcal{V}_2 = \mathcal{X}_1^{(1)}$ . This is the first fold;  $\mathcal{X}_i^{(j)}$  denotes half  $j$  of fold  $i$ .

To get the second fold, we shuffle  $\mathcal{X}$  randomly and divide this new fold into two,  $\mathcal{X}_2^{(1)}$  and  $\mathcal{X}_2^{(2)}$ . This can be implemented by drawing these from  $\mathcal{X}$  randomly without replacement, namely,  $\mathcal{X}_1^{(1)} \cup \mathcal{X}_1^{(2)} = \mathcal{X}_2^{(1)} \cup \mathcal{X}_2^{(2)} = \mathcal{X}$ . We then swap these two halves to get another pair. We do this for three more folds and because from each fold, we get two pairs, doing five folds, we get ten training and validation sets:

$$\begin{array}{ll} \mathcal{T}_1 = \mathcal{X}_1^{(1)} & \mathcal{V}_1 = \mathcal{X}_1^{(2)} \\ \mathcal{T}_2 = \mathcal{X}_1^{(2)} & \mathcal{V}_2 = \mathcal{X}_1^{(1)} \\ \mathcal{T}_3 = \mathcal{X}_2^{(1)} & \mathcal{V}_3 = \mathcal{X}_2^{(2)} \\ \mathcal{T}_4 = \mathcal{X}_2^{(2)} & \mathcal{V}_4 = \mathcal{X}_2^{(1)} \\ \vdots & \\ \mathcal{T}_9 = \mathcal{X}_5^{(1)} & \mathcal{V}_9 = \mathcal{X}_5^{(2)} \\ \mathcal{T}_{10} = \mathcal{X}_5^{(2)} & \mathcal{V}_{10} = \mathcal{X}_5^{(1)} \end{array}$$

Of course, we can do this for more than five folds and get more training/validation sets, but Dietterich (1998) points out that after five folds, the sets share many instances and overlap so much that the statistics calculated from these sets, namely, validation error rates, become too dependent and do not add new information. Even with five folds, the sets overlap and the statistics are dependent, but we can get away with this until five folds. On the other hand, if we do have fewer than five folds, we get less data (fewer than ten sets) and will not have a large enough sample to fit a distribution to and test our hypothesis on.

**Table 19.1** Confusion matrix for two classes.

True Class	Predicted class		Total
	Positive	Negative	
Positive	$tp$ : true positive	$fn$ : false negative	$p$
Negative	$fp$ : false positive	$tn$ : true negative	$n$
Total	$p'$	$n'$	$N$

### 19.6.3 Bootstrapping

BOOTSTRAP

To generate multiple samples from a single sample, an alternative to cross-validation is the *bootstrap* that generates new samples by drawing instances from the original sample *with* replacement. We saw the use of bootstrapping in section 17.6 to generate training sets for different learners in bagging. The bootstrap samples may overlap more than cross-validation samples and hence their estimates are more dependent; but is considered the best way to do resampling for very small datasets.

In the bootstrap, we sample  $N$  instances from a dataset of size  $N$  with replacement. The original dataset is used as the validation set. The probability that we pick an instance is  $1/N$ ; the probability that we do not pick it is  $1 - 1/N$ . The probability that we do not pick it after  $N$  draws is

$$\left(1 - \frac{1}{N}\right)^N \approx e^{-1} = 0.368$$

This means that the training data contains approximately 63.2 percent of the instances; that is, the system will not have been trained on 36.8 percent of the data, and the error estimate will be pessimistic. The solution is replication, that is, to repeat the process many times and look at the average behavior.

## 19.7 Measuring Classifier Performance

For classification, especially for two-class problems, a variety of measures has been proposed. There are four possible cases, as shown in table 19.1. For a positive example, if the prediction is also positive, this is a *true positive*; if our prediction is negative for a positive example, this is a *false negative*. For a negative example, if the prediction is also negative, we

**Table 19.2** Performance measures used in two-class problems.

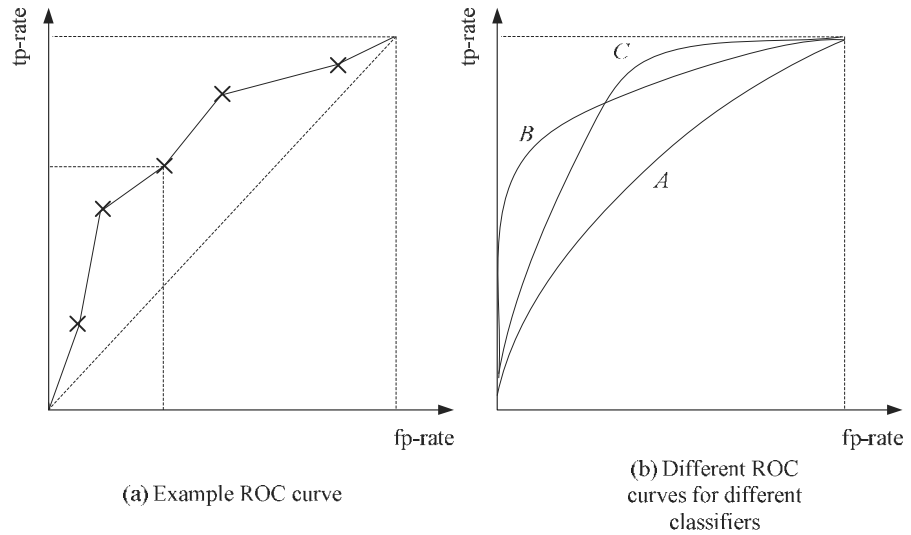
Name	Formula
error	$(fp + fn)/N$
accuracy	$(tp + tn)/N = 1 - \text{error}$
tp-rate	$tp/p$
fp-rate	$fp/n$
precision	$tp/p'$
recall	$tp/p = \text{tp-rate}$
sensitivity	$tp/p = \text{tp-rate}$
specificity	$tn/n = 1 - \text{fp-rate}$

have a *true negative*, and we have a *false positive* if we predict a negative example as positive.

In some two-class problems, we make a distinction between the two classes and hence the two type of errors, false positives and false negatives. Different measures appropriate in different settings are given in table 19.2. Let us envisage an authentication application where, for example, users log on to their accounts by voice. A false positive is wrongly logging on an impostor and a false negative is refusing a valid user. It is clear that the two type of errors are not equally bad; the former is much worse. True positive rate, *tp-rate*, also known as *hit rate*, measures what proportion of valid users we authenticate and false positive rate, *fp-rate*, also known as *false alarm rate*, is the proportion of impostors we wrongly accept.

Let us say the system returns  $\hat{P}(C_1|x)$ , the probability of the positive class, and for the negative class, we have  $\hat{P}(C_2|x) = 1 - \hat{P}(C_1|x)$ , and we choose “positive” if  $\hat{P}(C_1|x) > \theta$ . If  $\theta$  is close to 1, we hardly choose the positive class; that is, we will have no false positives but also few true positives. As we decrease  $\theta$  to increase the number of true positives, we risk introducing false positives.

For different values of  $\theta$ , we can get a number of pairs of (tp-rate, fp-rate) values and by connecting them we get the *receiver operating characteristics* (ROC) curve, as shown in figure 19.3a. Note that different values of  $\theta$  correspond to different loss matrices for the two types of error and the ROC curve can also be seen as the behavior of a classifier



**Figure 19.3** (a) Typical ROC curve. Each classifier has a threshold that allows us to move over this curve, and we decide on a point, based on the relative importance of hits versus false alarms, namely, true positives and false positives. The area below the ROC curve is called AUC. (b) A classifier is preferred if its ROC curve is closer to the upper-left corner (larger AUC). *B* and *C* are preferred over *A*; *B* and *C* are preferred under different loss matrices.

under different loss matrices (see exercise 1).

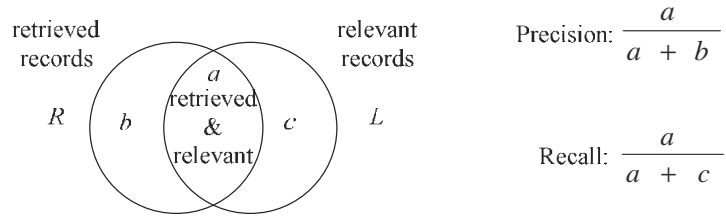
Ideally, a classifier has a tp-rate of 1 and a fp-rate of 0, and hence a classifier is better the more it gets closer to the upper-left corner. On the diagonal, we make as many true decisions as false ones, and this is the worst one can do (any classifier that is below the diagonal can be improved by flipping its decision). Given two classifiers, we can say one is better than the other one if it is above the other one; if two ROC curves intersect, we can say that the two classifiers are better under different loss conditions, as seen in figure 19.3b.

ROC allows a visual analysis; if we want to reduce the curve to a single number we can do this by calculating the *area under the curve* (AUC). A classifier ideally has an AUC of 1 and AUC values of different classifiers can be compared to give us a general performance averaged over different loss conditions.

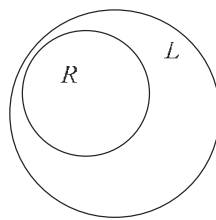
In *information retrieval*, there is a database of records; we make a

AREA UNDER THE  
CURVE

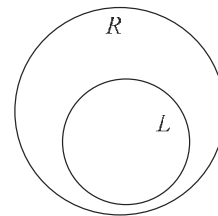
INFORMATION  
RETRIEVAL



(a) Precision and recall



(b) Precision = 1



(c) Recall = 1

**Figure 19.4** (a) Definition of precision and recall using Venn diagrams. (b) Precision is 1; all the retrieved records are relevant but there may be relevant ones not retrieved. (c) Recall is 1; all the relevant records are retrieved but there may also be irrelevant records that are retrieved.

PRECISION

RECALL

query, for example, by using some keywords, and a system (basically a two-class classifier) returns a number of records. In the database, there are relevant records and for a query, the system may retrieve some of them (true positives) but probably not all (false negatives); it may also wrongly retrieve records that are not relevant (false positives). The set of relevant and retrieved records can be visualized using a Venn diagram, as shown in figure 19.4a. *Precision* is the number of retrieved and relevant records divided by the total number of retrieved records; if precision is 1, all the retrieved records may be relevant but there may still be records that are relevant but not retrieved. *Recall* is the number of retrieved relevant records divided by the total number of relevant records; even if recall is 1, all the relevant records may be retrieved but there may also be irrelevant records that are retrieved, as shown in figure19.4c. As in the ROC curve, for different threshold values, one can draw a curve for precision vs. recall.

SENSITIVITY  
SPECIFICITY

From another perspective but with the same aim, there are the two measures of *sensitivity* and *specificity*. Sensitivity is the same as tp-rate and recall. Specificity is how well we detect the negatives, which is the number of true negatives divided by the total number of negatives; this is equal to 1 minus the false alarm rate. One can also draw a sensitivity vs. specificity curve using different thresholds.

CLASS CONFUSION  
MATRIX

In the case of  $K > 2$  classes, if we are using 0/1 error, the *class confusion matrix* is a  $K \times K$  matrix whose entry  $(i, j)$  contains the number of instances that belong to  $C_i$  but are assigned to  $C_j$ . Ideally, all off-diagonals should be 0, for no misclassification. The class confusion matrix allows us to pinpoint what types of misclassification occur, namely, if there are two classes that are frequently confused. Or, one can define  $K$  separate two-class problems, each one separating one class from the other  $K - 1$ .

### 19.8 Interval Estimation

INTERVAL ESTIMATION

Let us now do a quick review of *interval estimation* that we will use in hypothesis testing. A point estimator, for example, the maximum likelihood estimator, specifies a value for a parameter  $\theta$ . In interval estimation, we specify an interval within which  $\theta$  lies with a certain degree of confidence. To obtain such an interval estimator, we make use of the probability distribution of the point estimator.

For example, let us say we are trying to estimate the mean  $\mu$  of a normal density from a sample  $X = \{x^t\}_{t=1}^N$ .  $m = \sum_t x^t / N$  is the sample average and is the point estimator to the mean.  $m$  is the sum of normals and therefore is also normal,  $m \sim \mathcal{N}(\mu, \sigma^2 / N)$ . We define the statistic with a *unit normal distribution*:

UNIT NORMAL  
DISTRIBUTION

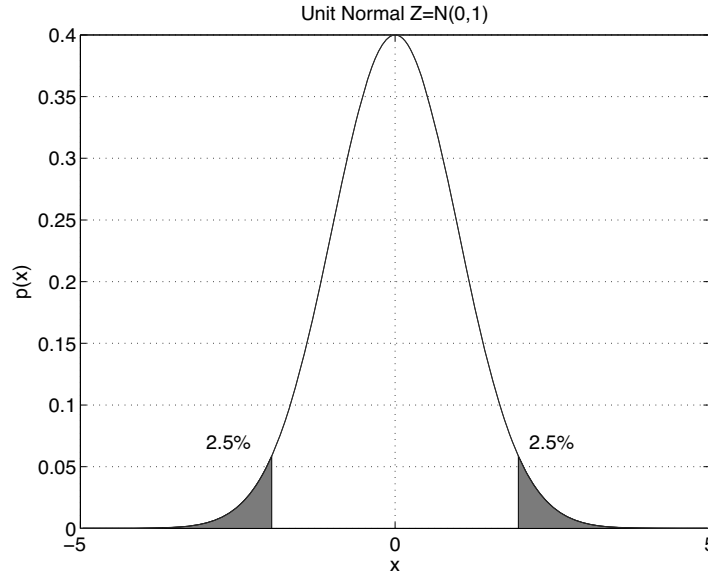
$$(19.1) \quad \frac{(m - \mu)}{\sigma / \sqrt{N}} \sim Z$$

We know that 95 percent of  $Z$  lies in  $(-1.96, 1.96)$ , namely,  $P\{-1.96 < Z < 1.96\} = 0.95$ , and we can write (see figure 19.5)

$$P \left\{ -1.96 < \sqrt{N} \frac{(m - \mu)}{\sigma} < 1.96 \right\} = 0.95$$

or equivalently

$$P \left\{ m - 1.96 \frac{\sigma}{\sqrt{N}} < \mu < m + 1.96 \frac{\sigma}{\sqrt{N}} \right\} = 0.95$$



**Figure 19.5** 95 percent of the unit normal distribution lies between  $-1.96$  and  $1.96$ .

TWO-SIDED  
CONFIDENCE  
INTERVAL

That is “with 95 percent confidence,”  $\mu$  will lie within  $1.96\sigma/\sqrt{N}$  units of the sample average. This is a *two-sided confidence interval*. With 99 percent confidence,  $\mu$  will lie in  $(m - 2.58\sigma/\sqrt{N}, m + 2.58\sigma/\sqrt{N})$ ; that is, if we want more confidence, the interval gets larger. The interval gets smaller as  $N$ , the sample size, increases.

This can be generalized for any required confidence as follows. Let us denote  $z_\alpha$  such that

$$P\{Z > z_\alpha\} = \alpha, \quad 0 < \alpha < 1$$

Because  $Z$  is symmetric around the mean,  $z_{1-\alpha/2} = -z_{\alpha/2}$ , and  $P\{X < -z_{\alpha/2}\} = P\{X > z_{\alpha/2}\} = \alpha/2$ . Hence for any specified level of confidence  $1 - \alpha$ , we have

$$P\{-z_{\alpha/2} < Z < z_{\alpha/2}\} = 1 - \alpha$$

and

$$P\left\{-z_{\alpha/2} < \sqrt{N} \frac{(m - \mu)}{\sigma} < z_{\alpha/2}\right\} = 1 - \alpha$$

or

$$(19.2) \quad P \left\{ m - z_{\alpha/2} \frac{\sigma}{\sqrt{N}} < \mu < m + z_{\alpha/2} \frac{\sigma}{\sqrt{N}} \right\} = 1 - \alpha$$

Hence a  $100(1 - \alpha)$  percent two-sided confidence interval for  $\mu$  can be computed for any  $\alpha$ .

Similarly, knowing that  $P\{Z < 1.64\} = 0.95$ , we have (see figure 19.6)

$$P \left\{ \sqrt{N} \frac{(m - \mu)}{\sigma} < 1.64 \right\} = 0.95$$

or

$$P \left\{ m - 1.64 \frac{\sigma}{\sqrt{N}} < \mu \right\} = 0.95$$

ONE-SIDED  
CONFIDENCE  
INTERVAL

and  $(m - 1.64\sigma/\sqrt{N}, \infty)$  is a 95 percent *one-sided upper confidence interval* for  $\mu$ , which defines a lower bound. Generalizing, a  $100(1 - \alpha)$  percent one-sided confidence interval for  $\mu$  can be computed from

$$(19.3) \quad P \left\{ m - z_{\alpha} \frac{\sigma}{\sqrt{N}} < \mu \right\} = 1 - \alpha$$

Similarly, the one-sided lower confidence interval that defines an upper bound can also be calculated.

In the previous intervals, we used  $\sigma$ ; that is, we assumed that the variance is known. If it is not, one can plug the sample variance

$$S^2 = \sum_t (x^t - m)^2 / (N - 1)$$

instead of  $\sigma^2$ . We know that when  $x^t \sim \mathcal{N}(\mu, \sigma^2)$ ,  $(N - 1)S^2/\sigma^2$  is chi-square with  $N - 1$  degrees of freedom. We also know that  $m$  and  $S^2$  are independent. Then,  $\sqrt{N}(m - \mu)/S$  is *t-distributed* with  $N - 1$  degrees of freedom (section A.3.7), denoted as

$$(19.4) \quad \frac{\sqrt{N}(m - \mu)}{S} \sim t_{N-1}$$

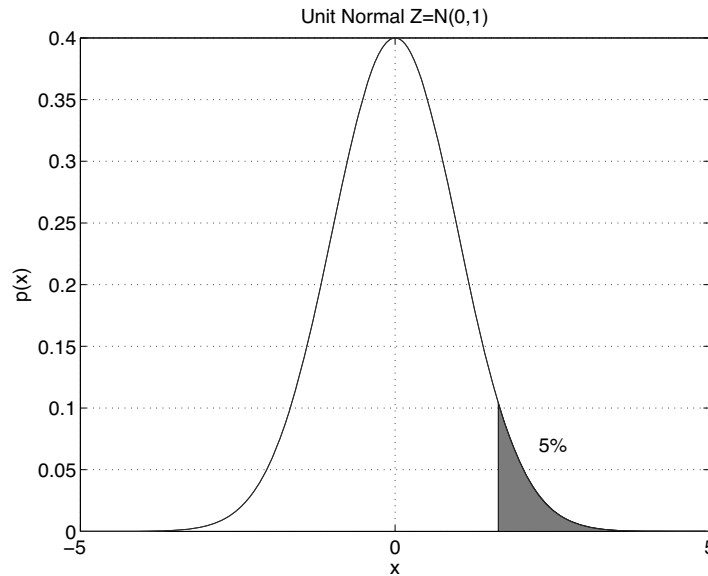
t DISTRIBUTION

Hence for any  $\alpha \in (0, 1/2)$ , we can define an interval, using the values specified by the *t distribution*, instead of the unit normal  $Z$

$$P \left\{ t_{1-\alpha/2, N-1} < \sqrt{N} \frac{(m - \mu)}{S} < t_{\alpha/2, N-1} \right\} = 1 - \alpha$$

or using  $t_{1-\alpha/2, N-1} = -t_{\alpha/2, N-1}$ , we can write

$$P \left\{ m - t_{\alpha/2, N-1} \frac{S}{\sqrt{N}} < \mu < m + t_{\alpha/2, N-1} \frac{S}{\sqrt{N}} \right\} = 1 - \alpha$$



**Figure 19.6** 95 percent of the unit normal distribution lies before 1.64.

Similarly, one-sided confidence intervals can be defined. The  $t$  distribution has larger spread (longer tails) than the unit normal distribution, and generally the interval given by the  $t$  is larger; this should be expected since additional uncertainty exists due to the unknown variance.

## 19.9 Hypothesis Testing

Instead of explicitly estimating some parameters, in certain applications we may want to use the sample to test some particular hypothesis concerning the parameters. For example, instead of estimating the mean, we may want to test whether the mean is less than 0.02. If the random sample is consistent with the hypothesis under consideration, we “fail to reject” the hypothesis; otherwise, we say that it is “rejected.” But when we make such a decision, we are not really saying that it is true or false but rather that the sample data appears to be consistent with it to a given degree of confidence or not.

HYPOTHESIS TESTING

In *hypothesis testing*, the approach is as follows. We define a statistic

**Table 19.3** Type I error, type II error, and power of a test.

	Decision	
Truth	Fail to reject	Reject
True	Correct	Type I error
False	Type II error	Correct (power)

that obeys a certain distribution if the hypothesis is correct. If the statistic calculated from the sample has very low probability of being drawn from this distribution, then we reject the hypothesis; otherwise, we fail to reject it.

Let us say we have a sample from a normal distribution with unknown mean  $\mu$  and known variance  $\sigma^2$ , and we want to test a specific hypothesis about  $\mu$ , for example, whether it is equal to a specified constant  $\mu_0$ . It is denoted as  $H_0$  and is called the *null hypothesis*

$$H_0 : \mu = \mu_0$$

against the alternative hypothesis

$$H_1 : \mu \neq \mu_0$$

$m$  is the point estimate of  $\mu$ , and it is reasonable to reject  $H_0$  if  $m$  is too far from  $\mu_0$ . This is where the interval estimate is used. We fail to reject the hypothesis with *level of significance*  $\alpha$  if  $\mu_0$  lies in the  $100(1 - \alpha)$  percent confidence interval, namely, if

$$(19.5) \quad \frac{\sqrt{N}(m - \mu_0)}{\sigma} \in (-z_{\alpha/2}, z_{\alpha/2})$$

We reject the null hypothesis if it falls outside, on either side. This is a *two-sided test*.

If we reject when the hypothesis is correct, this is a *type I error* and thus  $\alpha$ , set before the test, defines how much type I error we can tolerate, typical values being  $\alpha = 0.1, 0.05, 0.01$  (see table 19.3). A *type II error* is if we fail to reject the null hypothesis when the true mean  $\mu$  is unequal to  $\mu_0$ . The probability that  $H_0$  is not rejected when the true mean is  $\mu$  is a function of  $\mu$  and is given as

$$(19.6) \quad \beta(\mu) = P_{\mu} \left\{ -z_{\alpha/2} \leq \frac{m - \mu_0}{\sigma/\sqrt{N}} \leq z_{\alpha/2} \right\}$$

POWER FUNCTION  $1 - \beta(\mu)$  is called the *power function* of the test and is equal to the probability of rejection when  $\mu$  is the true value. Type II error probability increases as  $\mu$  and  $\mu_0$  gets closer, and we can calculate how large a sample we need for us to be able to detect a difference  $\delta = |\mu - \mu_0|$  with sufficient power.

ONE-SIDED TEST One can also have a *one-sided test* of the form

$$H_0 : \mu \leq \mu_0 \text{ vs } H_1 : \mu > \mu_0$$

as opposed to the two-sided test when the alternative hypothesis is  $\mu \neq \mu_0$ . The one-sided test with  $\alpha$  level of significance defines the  $100(1 - \alpha)$  confidence interval bounded on one side in which  $m$  should lie for the hypothesis not to be rejected. We fail to reject if

$$(19.7) \quad \frac{\sqrt{N}}{\sigma}(m - \mu_0) \in (-\infty, z_\alpha)$$

and reject outside. Note that the null hypothesis  $H_0$  also allows equality, which means that we get ordering information only if the test rejects. This tells us which of the two one-sided tests we should use. Whatever claim we have should be in  $H_1$  so that rejection of the test would support our claim.

If the variance is unknown, just as we did in the interval estimates, we use the sample variance instead of the population variance and the fact that

$$(19.8) \quad \frac{\sqrt{N}(m - \mu_0)}{S} \sim t_{N-1}$$

For example, for  $H_0 : \mu = \mu_0$  vs  $H_1 : \mu \neq \mu_0$ , we fail to reject at significance level  $\alpha$  if

$$(19.9) \quad \frac{\sqrt{N}(m - \mu_0)}{S} \in (-t_{\alpha/2, N-1}, t_{\alpha/2, N-1})$$

$t$  TEST which is known as the *two-sided t test*. A one-sided  $t$  test can be defined similarly.

## 19.10 Assessing a Classification Algorithm's Performance

Now that we have reviewed hypothesis testing, we are ready to see how it is used in testing error rates. We will discuss the case of classification error, but the same methodology applies for squared error in regression, log likelihoods in unsupervised learning, expected reward in

reinforcement learning, and so on, as long as we can write the appropriate parametric form for the sampling distribution. We will also discuss nonparametric tests when no such parametric form can be found.

We now start with error rate assessment, and, in the next section, we discuss error rate comparison.

### 19.10.1 Binomial Test

Let us start with the case where we have a single training set  $\mathcal{T}$  and a single validation set  $\mathcal{V}$ . We train our classifier on  $\mathcal{T}$  and test it on  $\mathcal{V}$ . We denote by  $p$  the probability that the classifier makes a misclassification error. We do not know  $p$ ; it is what we would like to estimate or test a hypothesis about. On the instance with index  $t$  from the validation set  $\mathcal{V}$ , let us say  $x^t$  denotes the correctness of the classifier's decision:  $x^t$  is a 0/1 Bernoulli random variable that takes the value 1 when the classifier commits an error and 0 when the classifier is correct. The binomial random variable  $X$  denotes the total number of errors:

$$X = \sum_{t=1}^N x^t$$

We would like to test whether the error probability  $p$  is less than or equal to some value  $p_0$  we specify:

$$H_0 : p \leq p_0 \text{ vs. } H_1 : p > p_0$$

If the probability of error is  $p$ , the probability that the classifier commits  $j$  errors out of  $N$  is

$$P\{X = j\} = \binom{N}{j} p^j (1-p)^{N-j}$$

BINOMIAL TEST It is reasonable to reject  $p \leq p_0$  if in such a case, the probability that we see  $X = e$  errors or more is very unlikely. That is, the *binomial test* rejects the hypothesis if

$$(19.10) \quad P\{X \geq e\} = \sum_{x=e}^N \binom{N}{x} p_0^x (1-p_0)^{N-x} < \alpha$$

where  $\alpha$  is the significance, for example, 0.05.

### 19.10.2 Approximate Normal Test

If  $p$  is the probability of error, our point estimate is  $\hat{p} = X/N$ . Then, it is reasonable to reject the null hypothesis if  $\hat{p}$  is much larger than  $p_0$ . How large is large enough is given by the sampling distribution of  $\hat{p}$  and the significance  $\alpha$ .

Because  $X$  is the sum of independent random variables from the same distribution, the central limit theorem states that for large  $N$ ,  $X/N$  is approximately normal with mean  $p_0$  and variance  $p_0(1 - p_0)$ . Then

$$(19.11) \quad \frac{X/N - p_0}{\sqrt{p_0(1 - p_0)}} \sim \mathcal{Z}$$

APPROXIMATE  
NORMAL TEST

where  $\sim$  denotes “approximately distributed.” Then, using equation 19.7, the *approximate normal test* rejects the null hypothesis if this value for  $X = e$  is greater than  $z_\alpha$ .  $z_{0.05}$  is 1.64. This approximation will work well as long as  $N$  is not too small and  $p$  is not very close to 0 or 1; as a rule of thumb, we require  $Np \geq 5$  and  $N(1 - p) \geq 5$ .

### 19.10.3 $t$ Test

The two tests we discussed earlier use a single validation set. If we run the algorithm  $K$  times, on  $K$  training/validation set pairs, we get  $K$  error percentages,  $p_i, i = 1, \dots, K$  on the  $K$  validation sets. Let  $x_i^t$  be 1 if the classifier trained on  $\mathcal{T}_i$  makes a misclassification error on instance  $t$  of  $\mathcal{V}_i$ ;  $x_i^t$  is 0 otherwise. Then

$$p_i = \frac{\sum_{t=1}^N x_i^t}{N}$$

Given that

$$m = \frac{\sum_{i=1}^K p_i}{K}, \quad S^2 = \frac{\sum_{i=1}^K (p_i - m)^2}{K - 1}$$

from equation 19.8, we know that we have

$$(19.12) \quad \frac{\sqrt{K}(m - p_0)}{S} \sim t_{K-1}$$

and the  $t$  test rejects the null hypothesis that the classification algorithm has  $p_0$  or less error percentage at significance level  $\alpha$  if this value is greater than  $t_{\alpha, K-1}$ . Typically,  $K$  is taken as 10 or 30.  $t_{0.05, 9} = 1.83$  and  $t_{0.05, 29} = 1.70$ .

## 19.11 Comparing Two Classification Algorithms

Given two learning algorithms, we want to compare and test whether they construct classifiers that have the same expected error rate.

### 19.11.1 McNemar's Test

CONTINGENCY TABLE

Given a training set and a validation set, we use two algorithms to train two classifiers on the training set and test them on the validation set and compute their errors. A *contingency table*, like the one shown here, is an array of natural numbers in matrix form representing counts, or frequencies:

$e_{00}$ : Number of examples misclassified by both	$e_{01}$ : Number of examples misclassified by 1 but not 2
$e_{10}$ : Number of examples misclassified by 2 but not 1	$e_{11}$ : Number of examples correctly classified by both

Under the null hypothesis that the classification algorithms have the same error rate, we expect  $e_{01} = e_{10}$  and these to be equal to  $(e_{01} + e_{10})/2$ . We have the chi-square statistic with one degree of freedom

$$(19.13) \quad \frac{(|e_{01} - e_{10}| - 1)^2}{e_{01} + e_{10}} \sim \chi_1^2$$

MCNEMAR'S TEST

and *McNemar's test* rejects the hypothesis that the two classification algorithms have the same error rate at significance level  $\alpha$  if this value is greater than  $\chi_{\alpha,1}^2$ . For  $\alpha = 0.05$ ,  $\chi_{0.05,1}^2 = 3.84$ .

### 19.11.2 K-Fold Cross-Validated Paired t Test

This set uses  $K$ -fold cross-validation to get  $K$  training/validation set pairs. We use the two classification algorithms to train on the training sets  $\mathcal{T}_i, i = 1, \dots, K$ , and test on the validation sets  $\mathcal{V}_i$ . The error percentages of the classifiers on the validation sets are recorded as  $p_i^1$  and  $p_i^2$ .

PAIRED TEST

If the two classification algorithms have the same error rate, then we expect them to have the same mean, or equivalently, that the difference of their means is 0. The difference in error rates on fold  $i$  is  $p_i = p_i^1 - p_i^2$ . This is a *paired test*; that is, for each  $i$ , both algorithms see the same training and validation sets. When this is done  $K$  times, we have a distribution of  $p_i$  containing  $K$  points. Given that  $p_i^1$  and  $p_i^2$  are both (approximately)

normal, their difference  $p_i$  is also normal. The null hypothesis is that this distribution has 0 mean:

$$H_0 : \mu = 0 \text{ vs. } H_1 : \mu \neq 0$$

We define

$$m = \frac{\sum_{i=1}^K p_i}{K}, \quad S^2 = \frac{\sum_{i=1}^K (p_i - m)^2}{K - 1}$$

Under the null hypothesis that  $\mu = 0$ , we have a statistic that is  $t$ -distributed with  $K - 1$  degrees of freedom:

$$(19.14) \quad \frac{\sqrt{K}(m - 0)}{S} = \frac{\sqrt{K} \cdot m}{S} \sim t_{K-1}$$

*K*-FOLD CV PAIRED *t*  
TEST

Thus the *K-fold cv paired t test* rejects the hypothesis that two classification algorithms have the same error rate at significance level  $\alpha$  if this value is outside the interval  $(-t_{\alpha/2, K-1}, t_{\alpha/2, K-1})$ .  $t_{0.025, 9} = 2.26$  and  $t_{0.025, 29} = 2.05$ .

If we want to test whether the first algorithm has less error than the second, we need a one-sided hypothesis and use a one-tailed test:

$$H_0 : \mu \geq 0 \text{ vs. } H_1 : \mu < 0$$

If the test rejects, our claim that the first one has significantly less error is supported.

### 19.11.3 $5 \times 2$ cv Paired *t* Test

In the  $5 \times 2$  cv *t* test, proposed by Dietterich (1998), we perform five replications of twofold cross-validation. In each replication, the dataset is divided into two equal-sized sets.  $p_i^{(j)}$  is the difference between the error rates of the two classifiers on fold  $j = 1, 2$  of replication  $i = 1, \dots, 5$ . The average on replication  $i$  is  $\bar{p}_i = (p_i^{(1)} + p_i^{(2)})/2$ , and the estimated variance is  $s_i^2 = (p_i^{(1)} - \bar{p}_i)^2 + (p_i^{(2)} - \bar{p}_i)^2$ .

Under the null hypothesis that the two classification algorithms have the same error rate,  $p_i^{(j)}$  is the difference of two identically distributed proportions, and ignoring the fact that these proportions are not independent,  $p_i^{(j)}$  can be treated as approximately normal distributed with 0 mean and unknown variance  $\sigma^2$ . Then  $p_i^{(j)}/\sigma$  is approximately unit normal. If we assume  $p_i^{(1)}$  and  $p_i^{(2)}$  are independent normals (which is not strictly true because their training and test sets are not drawn independently of each other), then  $s_i^2/\sigma^2$  has a chi-square distribution with

one degree of freedom. If each of the  $s_i^2$  are assumed to be independent (which is not true because they are all computed from the same set of available data), then their sum is chi-square with five degrees of freedom:

$$M = \frac{\sum_{i=1}^5 s_i^2}{\sigma^2} \sim \chi_5^2$$

and

$$(19.15) \quad t = \frac{p_1^{(1)}/\sigma}{\sqrt{M/5}} = \frac{p_1^{(1)}}{\sqrt{\sum_{i=1}^5 s_i^2/5}} \sim t_5$$

5 × 2 CV PAIRED t TEST

giving us a  $t$  statistic with five degrees of freedom. The  $5 \times 2$  *cv paired t test* rejects the hypothesis that the two classification algorithms have the same error rate at significance level  $\alpha$  if this value is outside the interval  $(-t_{\alpha/2,5}, t_{\alpha/2,5})$ .  $t_{0.025,5} = 2.57$ .

#### 19.11.4 5 × 2 cv Paired F Test

We note that the numerator in equation 19.15,  $p_1^{(1)}$ , is arbitrary; actually, ten different values can be placed in the numerator, namely,  $p_i^{(j)}$ ,  $j = 1, 2, i = 1, \dots, 5$ , leading to ten possible statistics:

$$(19.16) \quad t_i^{(j)} = \frac{p_i^{(j)}}{\sqrt{\sum_{i=1}^5 s_i^2/5}}$$

Alpaydm (1999) proposed an extension to the  $5 \times 2$  *cv t test* that combines the results of the ten possible statistics. If  $p_i^{(j)}/\sigma \sim \mathcal{Z}$ , then  $(p_i^{(j)})^2/\sigma^2 \sim \chi_1^2$  and their sum is chi-square with ten degrees of freedom:

$$N = \frac{\sum_{i=1}^5 \sum_{j=1}^2 (p_i^{(j)})^2}{\sigma^2} \sim \chi_{10}^2$$

Placing this in the numerator of equation 19.15, we get a statistic that is the ratio of two chi-square distributed random variables. Two such variables divided by their respective degrees of freedom is  $F$ -distributed with ten and five degrees of freedom (section A.3.8):

$$(19.17) \quad f = \frac{N/10}{M/5} = \frac{\sum_{i=1}^5 \sum_{j=1}^2 (p_i^{(j)})^2}{2 \sum_{i=1}^5 s_i^2} \sim F_{10,5}$$

5 × 2 CV PAIRED F TEST

$5 \times 2$  *cv paired F test* rejects the hypothesis that the classification algorithms have the same error rate at significance level  $\alpha$  if this value is greater than  $F_{\alpha,10,5}$ .  $F_{0.05,10,5} = 4.74$ .

## 19.12 Comparing Multiple Algorithms: Analysis of Variance

In many cases, we have more than two algorithms, and we would like to compare their expected error. Given  $L$  algorithms, we train them on  $K$  training sets, induce  $K$  classifiers with each algorithm, and then test them on  $K$  validation sets and record their error rates. This gives us  $L$  groups of  $K$  values. The problem then is the comparison of these  $L$  samples for statistically significant difference. This is an experiment with a single factor with  $L$  levels, the learning algorithms, and there are  $K$  replications for each level.

ANALYSIS OF  
VARIANCE

In *analysis of variance* (ANOVA), we consider  $L$  independent samples, each of size  $K$ , composed of normal random variables of unknown mean  $\mu_j$  and unknown common variance  $\sigma^2$ :

$$X_{ij} \sim \mathcal{N}(\mu_j, \sigma^2), j = 1, \dots, L, i = 1, \dots, K,$$

We are interested in testing the hypothesis  $H_0$  that all means are equal:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_L \text{ vs. } H_1 : \mu_r \neq \mu_s, \text{ for at least one pair } (r, s)$$

The comparison of error rates of multiple classification algorithms fits this scheme. We have  $L$  classification algorithms, and we have their error rates on  $K$  validation folds.  $X_{ij}$  is the number of validation errors made by the classifier, which is trained by classification algorithm  $j$  on fold  $i$ . Each  $X_{ij}$  is binomial and approximately normal. If  $H_0$  is not rejected, we fail to find a significant error difference among the error rates of the  $L$  classification algorithms. This is therefore a generalization of the tests we saw in section 19.11 that compared the error rates of two classification algorithms. The  $L$  classification algorithms may be different or may use different hyperparameters, for example, number of hidden units in a multilayer perceptron, number of neighbors in  $k$ -nn, and so forth.

The approach in ANOVA is to derive two estimators of  $\sigma^2$ . One estimator is designed such that it is true only when  $H_0$  is true, and the second is always a valid estimator, regardless of whether  $H_0$  is true or not. ANOVA then rejects  $H_0$ , namely, that the  $L$  samples are drawn from the same population, if the two estimators differ significantly.

Our first estimator to  $\sigma^2$  is valid only if the hypothesis is true, namely,  $\mu_j = \mu, j = 1, \dots, L$ . If  $X_{ij} \sim \mathcal{N}(\mu, \sigma^2)$ , then the group average

$$m_j = \sum_{i=1}^K \frac{X_{ij}}{K}$$

is also normal with mean  $\mu$  and variance  $\sigma^2/K$ . If the hypothesis is true, then  $m_j, j = 1, \dots, L$  are  $L$  instances drawn from  $\mathcal{N}(\mu, \sigma^2/K)$ . Then *their* mean and variance are

$$m = \frac{\sum_{j=1}^L m_j}{L}, \quad S^2 = \frac{\sum_j (m_j - m)^2}{L-1}$$

Thus an estimator of  $\sigma^2$  is  $K \cdot S^2$ , namely,

$$(19.18) \quad \hat{\sigma}_b^2 = K \sum_{j=1}^L \frac{(m_j - m)^2}{L-1}$$

Each of  $m_j$  is normal and  $(L-1)S^2/(\sigma^2/K)$  is chi-square with  $(L-1)$  degrees of freedom. Then, we have

$$(19.19) \quad \sum_j \frac{(m_j - m)^2}{\sigma^2/K} \sim \chi_{L-1}^2$$

We define  $SS_b$ , the between-group sum of squares, as

$$SS_b \equiv K \sum_j (m_j - m)^2$$

So, when  $H_0$  is true, we have

$$(19.20) \quad \frac{SS_b}{\sigma^2} \sim \chi_{L-1}^2$$

Our second estimator of  $\sigma^2$  is the average of group variances,  $S_j^2$ , defined as

$$S_j^2 = \frac{\sum_{i=1}^K (X_{ij} - m_j)^2}{K-1}$$

and their average is

$$(19.21) \quad \hat{\sigma}_w^2 = \sum_{j=1}^L \frac{S_j^2}{L} = \sum_j \sum_i \frac{(X_{ij} - m_j)^2}{L(K-1)}$$

We define  $SS_w$ , the within-group sum of squares:

$$SS_w \equiv \sum_j \sum_i (X_{ij} - m_j)^2$$

Remembering that for a normal sample, we have

$$(K-1) \frac{S_j^2}{\sigma^2} \sim \chi_{K-1}^2$$

and that the sum of chi-squares is also a chi-square, we have

$$(K-1) \sum_{j=1}^L \frac{S_j^2}{\sigma^2} \sim \chi_{L(K-1)}^2$$

So

$$(19.22) \quad \frac{SS_w}{\sigma^2} \sim \chi_{L(K-1)}^2$$

Then we have the task of comparing two variances for equality, which we can do by checking whether their ratio is close to 1. The ratio of two independent chi-square random variables divided by their respective degrees of freedom is a random variable that is  $F$ -distributed, and hence when  $H_0$  is true, we have

$$(19.23) \quad F_0 = \left( \frac{SS_b/\sigma^2}{L-1} \right) \bigg/ \left( \frac{SS_w/\sigma^2}{L(K-1)} \right) = \frac{SS_b/(L-1)}{SS_w/(L(K-1))} = \frac{\hat{\sigma}_b^2}{\hat{\sigma}_w^2} \sim F_{L-1, L(K-1)}$$

For any given significance value  $\alpha$ , the hypothesis that the  $L$  classification algorithms have the same expected error rate is rejected if this statistic is greater than  $F_{\alpha, L-1, L(K-1)}$ .

Note that we are rejecting if the two estimators disagree significantly. If  $H_0$  is not true, then the variance of  $m_j$  around  $m$  will be larger than what we would normally have if  $H_0$  were true, and hence if  $H_0$  is not true, the first estimator  $\hat{\sigma}_b^2$  will overestimate  $\sigma^2$ , and the ratio will be greater than 1. For  $\alpha = 0.05$ ,  $L = 5$  and  $K = 10$ ,  $F_{0.05, 4, 45} = 2.6$ . If  $X_{ij}$  vary around  $m$  with a variance of  $\sigma^2$ , then if  $H_0$  is true,  $m_j$  vary around  $m$  by  $\sigma^2/K$ . If it seems as if they vary more, then  $H_0$  should be rejected because the displacement of  $m_j$  around  $m$  is more than what can be explained by some constant added noise.

The name *analysis of variance* is derived from a partitioning of the total variability in the data into its components.

$$(19.24) \quad SS_T \equiv \sum_j \sum_i (X_{ij} - m)^2$$

$SS_T$  divided by its degree of freedom, namely,  $K \cdot L - 1$  (there are  $K \cdot L$  data points, and we lose one degree of freedom because  $m$  is fixed), gives us the sample variance of  $X_{ij}$ . It can be shown that (exercise 5) the total sum of squares can be split into between-group sum of squares and within-group sum of squares

$$(19.25) \quad SS_T = SS_b + SS_w$$

**Table 19.4** The analysis of variance (ANOVA) table for a single factor model.

Source of variation	Sum of squares	Degrees of freedom	Mean square	$F_0$
Between groups	$SS_b \equiv K \sum_j (m_j - m)^2$	$L - 1$	$MS_b = \frac{SS_b}{L-1}$	$\frac{MS_b}{MS_w}$
Within groups	$SS_w \equiv \sum_j \sum_i (X_{ij} - m_j)^2$	$L(K - 1)$	$MS_w = \frac{SS_w}{L(K-1)}$	
Total	$SS_T \equiv \sum_j \sum_i (X_{ij} - m)^2$	$L \cdot K - 1$		

Results of ANOVA are reported in an ANOVA table as shown in table 19.4. This is the basic *one-way* analysis of variance where there is a single factor, for example, learning algorithm. We may consider experiments with multiple factors, for example, we can have one factor for classification algorithms and another factor for feature extraction algorithms used before, and this will be a *two-factor experiment with interaction*.

POSTHOC TESTING

If the hypothesis is rejected, we only know that there is some difference between the  $L$  groups but we do not know where. For this, we do *posthoc testing*, that is, an additional set of tests involving subsets of groups, for example, pairs.

LEAST SQUARE DIFFERENCE TEST

Fisher's *least square difference test* (LSD) compares groups in a pairwise manner. For each group, we have  $m_i \sim \mathcal{N}(\mu_i, \sigma_w^2 = MS_w/K)$  and  $m_i - m_j \sim \mathcal{N}(\mu_i - \mu_j, 2\sigma_w^2)$ . Then, under the null hypothesis that  $H_0 : \mu_i = \mu_j$ , we have

$$t = \frac{m_i - m_j}{\sqrt{2}\sigma_w} \sim t_{L(K-1)}$$

We reject  $H_0$  in favor of the alternative hypothesis  $H_1 : \mu_1 \neq \mu_2$  if  $|t| > t_{\alpha/2, L(K-1)}$ . Similarly, one-sided tests can be defined to find pairwise orderings.

MULTIPLE COMPARISONS

When we do a number of tests to draw one conclusion, this is called *multiple comparisons*, and we need to keep in mind that if  $T$  hypotheses are to be tested, each at significance level  $\alpha$ , then the probability that at least one hypothesis is incorrectly rejected is at most  $T\alpha$ . For example,

the probability that six confidence intervals, each calculated at 95 percent individual confidence intervals, will simultaneously be correct is at least 70 percent. Thus to ensure that the overall confidence interval is at least  $100(1 - \alpha)$ , each confidence interval should be set at  $100(1 - \alpha/T)$ . This is called a *Bonferroni correction*.

BONFERRONI  
CORRECTION

Sometimes it may be the case that ANOVA rejects and none of the posthoc pairwise tests find a significant difference. In such a case, our conclusion is that there is a difference between the means but that we need more data to be able to pinpoint the source of the difference.

Note that the main cost is the training and testing of  $L$  classification algorithms on  $K$  training/validation sets. Once this is done and the values are stored in a  $K \times L$  table, calculating the ANOVA or pairwise comparison test statistics from those is very cheap in comparison.

### 19.13 Comparison over Multiple Datasets

Let us say we want to compare two or more algorithms on several datasets and not one. What makes this different is that an algorithm depending on how well its inductive bias matches the problem will behave differently on different datasets, and these error values on different datasets cannot be said to be normally distributed around some mean accuracy. This implies that the parametric tests that we discussed in the previous sections based on binomials being approximately normal are no longer applicable and we need to resort to *nonparametric tests*. The advantage of having such tests is that we can also use them for comparing other statistics that are not normal, for example, training times, number of free parameters, and so on.

NONPARAMETRIC  
TESTS

Parametric tests are generally robust to slight departures from normality, especially if the sample is large. Nonparametric tests are distribution free but are less efficient; that is, if both are applicable, a parametric test should be preferred. The corresponding nonparametric test will require a larger sample to achieve the same power. Nonparametric tests assume no knowledge about the distribution of the underlying population but only that the values can be compared or ordered, and, as we will see, such tests make use of this order information.

When we have an algorithm trained on a number of different datasets, the average of its errors on these datasets is not a meaningful value, and, for example, we cannot use such averages to compare two algorithms  $A$

and  $B$ . To compare two algorithms, the only piece of information we can use is if on any dataset,  $A$  is more accurate than  $B$ ; we can then count the number of times  $A$  is more accurate than  $B$  and check whether this could have been by chance if they indeed were equally accurate. With more than two algorithms, we will look at the average *ranks* of the learners trained by different algorithms. Nonparametric tests basically use this rank data and not the absolute values.

Before proceeding with the details of these tests, it should be stressed that it does not make sense to compare error rates of algorithms on a whole variety of applications. Because there is no such thing as the “best learning algorithm,” such tests would not be conclusive. However, we can compare algorithms on a number of datasets, or versions, of the same application. For example, we may have a number of different datasets for face recognition but with different properties (resolution, lighting, number of subjects, and so on), and we may use a nonparametric test to compare algorithms on those; different properties of the datasets would make it impossible for us to lump images from different datasets together in a single set, but we can train algorithms separately on different datasets, obtain ranks separately, and then combine these to get an overall decision.

### 19.13.1 Comparing Two Algorithms

Let us say we want to compare two algorithms. We both train and validate them on  $i = 1, \dots, N$  different datasets in a paired manner—that is, all the conditions except the different algorithms should be identical. We get results  $e_i^1$  and  $e_i^2$  and if we use  $K$ -fold cross-validation on each dataset, these are averages or medians of the  $K$  values. The *sign test* is based on the idea that if the two algorithms have equal error, on each dataset, there should be  $1/2$  probability that the first has less error than the second, and thus we expect the first to win on  $N/2$  datasets. Let us define

$$X_i = \begin{cases} 1 & \text{if } e_i^1 < e_i^2 \\ 0 & \text{otherwise} \end{cases} \quad \text{and } X = \sum_{i=1}^N X_i$$

Let us say we want to test

$$H_0 : \mu_1 \geq \mu_2 \text{ vs. } H_1 : \mu_1 < \mu_2$$

If the null hypothesis is correct,  $X$  is binomial in  $N$  trials with  $p = 1/2$ . Let us say that we saw that the first one wins on  $X = e$  datasets. Then, the probability that we have  $e$  or less wins when indeed  $p = 1/2$  is

$$P\{X \leq e\} = \sum_{x=0}^e \binom{N}{x} \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{N-x}$$

and we reject if this probability is too small, that is, less than  $\alpha$ . If there are ties, we divide them equally to both sides; that is, if there are  $t$  ties, we add  $t/2$  to  $e$  (if  $t$  is odd, we ignore the odd one and decrease  $N$  by 1).

In testing

$$H_0 : \mu_1 \leq \mu_2 \text{ vs. } H_1 : \mu_1 > \mu_2$$

we reject if  $P\{X \geq e\} < \alpha$ .

For the two-sided test

$$H_0 : \mu_1 = \mu_2 \text{ vs. } H_1 : \mu_1 \neq \mu_2$$

we reject if  $e$  is too small or too large. If  $e < N/2$ , we reject if  $2P\{X \leq e\} < \alpha$ ; if  $e > N/2$ , we reject if  $2P\{X \geq e\} < \alpha$ —we need to find the corresponding tail, and we multiply it by 2 because it is a two-tailed test.

As we discussed before, nonparametric tests can be used to compare any measurements, for example, training times. In such a case, we see the advantage of a nonparametric test that uses order rather than averages of absolute values. Let us say we compare two algorithms on ten datasets, nine of which are small and have training times for both algorithms on the order of minutes, and one that is very large and whose training time is on the order of a day. If we use a parametric test and take the average of training times, the single large dataset will dominate the decision, but when we use the nonparameric test and compare values separately on each dataset, using the order will have the effect of normalizing separately for each dataset and hence will help us make a robust decision.

We can also use the sign test as a one sample test, for example, to check if the average error on all datasets is less than two percent, by comparing  $\mu_1$  not by the mean of a second population but by a constant  $\mu_0$ . We can do this simply by plugging the constant  $\mu_0$  in place of all observations from a second sample and using the procedure used earlier; that is, we will count how many times we get more or less than 0.02 and check if this is too unlikely under the null hypothesis. For large  $N$ , normal

approximation to the binomial can be used (exercise 6), but in practice, the number of datasets may be smaller than 20. Note that the sign test is a test on the median of a population, which is equal to the mean if the distribution is symmetric.

WILCOXON SIGNED  
RANK TEST

The sign test only uses the sign of the difference and not its magnitude, but we may envisage a case where the first algorithm, when it wins, always wins by a large margin whereas the second algorithm, when it wins, always wins barely. The *Wilcoxon signed rank test* uses both the sign and the magnitude of differences, as follows:

Let us say, additional to the sign of differences, we also calculate  $m_i = |e_i^1 - e_i^2|$  and then we order them so that the smallest,  $\min_i m_i$ , is assigned rank 1, the next smallest is assigned rank 2, and so on. If there are ties, their ranks are given the average value that they would receive if they differed slightly. For example, if the magnitudes are 2, 1, 2, 4, the ranks are 2.5, 1, 2.5, 4. We then calculate  $w_+$  as the sum of all ranks whose signs are positive and  $w_-$  as the sum of all ranks whose signs are negative.

The null hypothesis  $\mu_1 \leq \mu_2$  can be rejected in favor of the alternative  $\mu_1 > \mu_2$  only if  $w_+$  is much smaller than  $w_-$ . Similarly, the two-sided hypothesis  $\mu_1 = \mu_2$  can be rejected in favor of the alternative  $\mu_1 \neq \mu_2$  only if either  $w_+$  or  $w_-$ , that is,  $w = \min(w_+, w_-)$ , is very small. The critical values for the Wilcoxon signed rank test are tabulated and for  $N > 20$ , normal approximations can be used.

### 19.13.2 Multiple Algorithms

KRUSKAL-WALLIS TEST

The *Kruskal-Wallis test* is the nonparametric version of ANOVA and is a multiple sample generalization of a rank test. Given the  $M = L \cdot N$  observations, for example, error rates, of  $L$  algorithms on  $N$  datasets,  $X_{ij}, i = 1, \dots, L, j = 1, \dots, N$ , we rank them from the smallest to the largest and assign them ranks,  $R_{ij}$ , between 1 and  $M$ , again taking averages in case of ties. If the null hypothesis

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_L$$

is true, then the average of ranks of algorithm  $i$  should be approximately halfway between 1 and  $M$ , that is,  $(M + 1)/2$ . We denote the sample average rank of algorithm  $i$  by  $\bar{R}_{i\cdot}$  and we reject the hypothesis if the average ranks seem to differ from halfway. The test statistic

$$H = \frac{12}{(M + 1)L} \sum_{i=1}^L \left( \bar{R}_{i\cdot} - \frac{M + 1}{2} \right)^2$$

is approximately chi-square distributed with  $L - 1$  degrees of freedom and we reject the null hypothesis if the statistic exceeds  $\chi_{\alpha, L-1}$ .

TUKEY'S TEST Just like the parametric ANOVA, if the null hypothesis is rejected, we can do posthoc testing to check for pairwise comparison of ranks. One method for this is *Tukey's test*, which makes use of the *studentized range statistic*

$$q = \frac{\bar{R}_{max} - \bar{R}_{min}}{\sigma_w}$$

where  $\bar{R}_{max}$  and  $\bar{R}_{min}$  are the largest and smallest means (of ranks), respectively, out of the  $L$  means, and  $\sigma_w^2$  is the average variance of ranks around group rank averages. We reject that groups  $i$  and  $j$  have the same ranks in favor of the alternative hypothesis that they are different if

$$|\bar{R}_{i\bullet} - \bar{R}_{j\bullet}| > q_{\alpha}(L, L(K - 1))\sigma_w$$

where  $q_{\alpha}(L, L(K - 1))$  are tabulated. One-sided tests can also be defined to order algorithms in terms of average rank.

Demsar (2006) proposes to use CD (critical difference) diagrams for visualization. On a scale of 1 to  $L$ , we mark the averages,  $\bar{R}_{i\bullet}$ , and draw lines of length given by the critical difference,  $q_{\alpha}(L, L(K - 1))\sigma_w$ , between groups, so that lines connect groups that are not statistically significantly different.

## 19.14 Notes

The material related to experiment design follows the discussion from (Montgomery 2005), which here is adapted for machine learning. A more detailed discussion of interval estimation, hypothesis testing, and analysis of variance can be found in any introductory statistics book, for example, Ross 1987.

Dietterich (1998) discusses statistical tests and compares them on a number of applications using different classification algorithms. A review of ROC use and AUC calculation is given in Fawcett 2006. Demsar (2006) reviews statistical tests for comparing classifiers over multiple datasets.

When we compare two or more algorithms, if the null hypothesis that they have the same error rate is not rejected, we choose the simpler one, namely, the one with less space or time complexity. That is, we use our prior preference if the data does not prefer one in terms of error rate. For example, if we compare a linear model and a nonlinear model and

if the test does not reject that they have the same expected error rate, we should go for the simpler linear model. Even if the test rejects, in choosing one algorithm over another, error rate is only one of the criteria. Other criteria like training (space/time) complexity, testing complexity, and interpretability may override in practical applications.

This is how the posthoc test results are used in the MultiTest algorithm (Yıldız and Alpaydın 2006) to generate a full ordering. We do  $L(L - 1)/2$  one-sided pairwise tests to order the  $L$  algorithms, but it is very likely that the tests will not give a full ordering but only a partial order. The missing links are filled in using the prior complexity information to get a full order. A topological sort gives an ordering of algorithms using both types of information, error and complexity.

There are also tests to allow checking for *contrasts*. Let us say 1 and 2 are neural network methods and 3 and 4 are fuzzy logic methods. We can then test whether the average of 1 and 2 differs from the average of 3 and 4, thereby allowing us to compare methods in general.

Another important point to note is that we are only assessing or comparing misclassifications. This implies that from our point of view, all misclassifications have the same cost. When this is not the case, our tests should be based on risks taking a suitable loss function into account. Not much work has been done in this area. Similarly, these tests should be generalized from classification to regression, so as to be able to assess the mean square errors of regression algorithms, or to be able to compare the errors of two regression algorithms.

In comparing two classification algorithms, note that we are testing only whether they have the same expected error rate. If they do, this does not mean that they make the same errors. This is an idea that we used in chapter 17; we can combine multiple models to improve accuracy if different classifiers make different errors.

## 19.15 Exercises

1. In a two-class problem, let us say we have the loss matrix where  $\lambda_{11} = \lambda_{22} = 0$ ,  $\lambda_{21} = 1$  and  $\lambda_{12} = \alpha$ . Determine the threshold of decision as a function of  $\alpha$ .
2. We can simulate a classifier with error probability  $p$  by drawing samples from a Bernoulli distribution. Doing this, implement the binomial, approximate, and  $t$  tests for  $p_0 \in (0, 1)$ . Repeat these tests at least 1,000 times for several values of  $p$  and calculate the probability of rejecting the null hypothesis. What do you expect the probability of reject to be when  $p_0 = p$ ?

3. Assume  $x^t \sim \mathcal{N}(\mu, \sigma^2)$  where  $\sigma^2$  is known. How can we test for  $H_0 : \mu \geq \mu_0$  vs.  $H_1 : \mu < \mu_0$ ?
4. The  $K$ -fold cross-validated  $t$  test only tests for the equality of error rates. If the test rejects, we do not know which classification algorithm has the lower error rate. How can we test whether the first classification algorithm does not have higher error rate than the second one? Hint: We have to test  $H_0 : \mu \leq 0$  vs.  $H_1 : \mu > 0$ .
5. Show that the total sum of squares can be split into between-group sum of squares and within-group sum of squares as  $SS_T = SS_b + SS_w$ .
6. Use the normal approximation to the binomial for the sign test.
7. Let us say we have three classification algorithms. How can we order these three from best to worst?
8. If we have two variants of algorithm  $A$  and three variants of algorithm  $B$ , how can we compare the overall accuracies of  $A$  and  $B$  taking all their variants into account?
9. Propose a suitable test to compare the errors of two regression algorithms.
10. Propose a suitable test to compare the expected rewards of two reinforcement learning algorithms.

## 19.16 References

- Alpaydm, E. 1999. "Combined  $5 \times 2$  cv  $F$  Test for Comparing Supervised Classification Learning Algorithms." *Neural Computation* 11: 1885–1892.
- Bouckaert, R. R. 2003. "Choosing between Two Learning Algorithms based on Calibrated Tests." In *Twentieth International Conference on Machine Learning*, ed. T. Fawcett and N. Mishra, 51–58. Menlo Park, CA: AAAI Press.
- Demsar, J. 2006. "Statistical Comparison of Classifiers over Multiple Data Sets." *Journal of Machine Learning Research* 7: 1–30.
- Dietterich, T. G. 1998. "Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms." *Neural Computation* 10: 1895–1923.
- Fawcett, T. 2006. "An Introduction to ROC Analysis." *Pattern Recognition Letters* 27: 861–874.
- Montgomery, D. C. 2005. *Design and Analysis of Experiments*. 6th ed., New York: Wiley.
- Ross, S. M. 1987. *Introduction to Probability and Statistics for Engineers and Scientists*. New York: Wiley.

- Turney, P. 2000. "Types of Cost in Inductive Concept Learning." Paper presented at Workshop on Cost-Sensitive Learning at the Seventeenth International Conference on Machine Learning, Stanford University, Stanford, CA, July 2.
- Wolpert, D. H. 1995. "The Relationship between PAC, the Statistical Physics Framework, the Bayesian Framework, and the VC Framework." In *The Mathematics of Generalization*, ed. D. H. Wolpert, 117-214. Reading, MA: Addison-Wesley.
- Yıldız, O. T., and E. Alpaydın. 2006. "Ordering and Finding the Best of  $K > 2$  Supervised Learning Algorithms." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28: 392-402.



# A *Probability*

*We review briefly the elements of probability, the concept of a random variable, and example distributions.*

## A.1 Elements of Probability

A RANDOM experiment is one whose outcome is not predictable with certainty in advance (Ross 1987; Casella and Berger 1990). The set of all possible outcomes is known as the *sample space*  $S$ . A sample space is *discrete* if it consists of a finite (or countably infinite) set of outcomes; otherwise it is *continuous*. Any subset  $E$  of  $S$  is an *event*. Events are sets, and we can talk about their complement, intersection, union, and so forth.

One interpretation of probability is as a *frequency*. When an experiment is continually repeated under the exact same conditions, for any event  $E$ , the proportion of time that the outcome is in  $E$  approaches some constant value. This constant limiting frequency is the probability of the event, and we denote it as  $P(E)$ .

Probability sometimes is interpreted as a *degree of belief*. For example, when we speak of Turkey's probability of winning the World Soccer Cup in 2010, we do not mean a frequency of occurrence, since the championship will happen only once and it has not yet occurred (at the time of the writing of this book). What we mean in such a case is a subjective degree of belief in the occurrence of the event. Because it is subjective, different individuals may assign different probabilities to the same event.

### A.1.1 Axioms of Probability

Axioms ensure that the probabilities assigned in a random experiment can be interpreted as relative frequencies and that the assignments are consistent with our intuitive understanding of relationships among relative frequencies:

1.  $0 \leq P(E) \leq 1$ . If  $E_1$  is an event that cannot possibly occur, then  $P(E_1) = 0$ . If  $E_2$  is sure to occur,  $P(E_2) = 1$ .
2.  $S$  is the sample space containing all possible outcomes,  $P(S) = 1$ .
3. If  $E_i, i = 1, \dots, n$  are mutually exclusive (i.e., if they cannot occur at the same time, as in  $E_i \cap E_j = \emptyset, j \neq i$ , where  $\emptyset$  is the *null event* that does not contain any possible outcomes), we have

$$(A.1) \quad P\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n P(E_i)$$

For example, letting  $E^c$  denote the *complement* of  $E$ , consisting of all possible outcomes in  $S$  that are not in  $E$ , we have  $E \cap E^c = \emptyset$  and

$$\begin{aligned} P(E \cup E^c) &= P(E) + P(E^c) = 1 \\ P(E^c) &= 1 - P(E) \end{aligned}$$

If the intersection of  $E$  and  $F$  is not empty, we have

$$(A.2) \quad P(E \cup F) = P(E) + P(F) - P(E \cap F)$$

### A.1.2 Conditional Probability

$P(E|F)$  is the probability of the occurrence of event  $E$  given that  $F$  occurred and is given as

$$(A.3) \quad P(E|F) = \frac{P(E \cap F)}{P(F)}$$

Knowing that  $F$  occurred reduces the sample space to  $F$ , and the part of it where  $E$  also occurred is  $E \cap F$ . Note that equation A.3 is well-defined only if  $P(F) > 0$ . Because  $\cap$  is commutative, we have

$$P(E \cap F) = P(E|F)P(F) = P(F|E)P(E)$$

which gives us *Bayes' formula*:

$$(A.4) \quad P(F|E) = \frac{P(E|F)P(F)}{P(E)}$$

When  $F_i$  are mutually exclusive and exhaustive, namely,  $\bigcup_{i=1}^n F_i = S$

$$(A.5) \quad \begin{aligned} E &= \bigcup_{i=1}^n E \cap F_i \\ P(E) &= \sum_{i=1}^n P(E \cap F_i) = \sum_{i=1}^n P(E|F_i)P(F_i) \end{aligned}$$

Bayes' formula allows us to write

$$(A.6) \quad P(F_i|E) = \frac{P(E \cap F_i)}{P(E)} = \frac{P(E|F_i)P(F_i)}{\sum_j P(E|F_j)P(F_j)}$$

If  $E$  and  $F$  are *independent*, we have  $P(E|F) = P(E)$  and thus

$$(A.7) \quad P(E \cap F) = P(E)P(F)$$

That is, knowledge of whether  $F$  has occurred does not change the probability that  $E$  occurs.

## A.2 Random Variables

A *random variable* is a function that assigns a number to each outcome in the sample space of a random experiment.

### A.2.1 Probability Distribution and Density Functions

The *probability distribution function*  $F(\cdot)$  of a random variable  $X$  for any real number  $a$  is

$$(A.8) \quad F(a) = P\{X \leq a\}$$

and we have

$$(A.9) \quad P\{a < X \leq b\} = F(b) - F(a)$$

If  $X$  is a discrete random variable

$$(A.10) \quad F(a) = \sum_{\forall x \leq a} P(x)$$

where  $P(\cdot)$  is the *probability mass function* defined as  $P(a) = P\{X = a\}$ . If  $X$  is a *continuous* random variable,  $p(\cdot)$  is the *probability density function* such that

$$(A.11) \quad F(a) = \int_{-\infty}^a p(x) dx$$

### A.2.2 Joint Distribution and Density Functions

In certain experiments, we may be interested in the relationship between two or more random variables, and we use the *joint* probability distribution and density functions of  $X$  and  $Y$  satisfying

$$(A.12) \quad F(x, y) = P\{X \leq x, Y \leq y\}$$

Individual *marginal* distributions and densities can be computed by marginalizing, namely, summing over the free variable:

$$(A.13) \quad F_X(x) = P\{X \leq x\} = P\{X \leq x, Y \leq \infty\} = F(x, \infty)$$

In the discrete case, we write

$$(A.14) \quad P(X = x) = \sum_j P(x, y_j)$$

and in the continuous case, we have

$$(A.15) \quad p_X(x) = \int_{-\infty}^{\infty} p(x, y) dy$$

If  $X$  and  $Y$  are *independent*, we have

$$(A.16) \quad p(x, y) = p_X(x)p_Y(y)$$

These can be generalized in a straightforward manner to more than two random variables.

### A.2.3 Conditional Distributions

When  $X$  and  $Y$  are random variables

$$(A.17) \quad P_{X|Y}(x|y) = P\{X = x|Y = y\} = \frac{P\{X = x, Y = y\}}{P\{Y = y\}} = \frac{P(x, y)}{P_Y(y)}$$

### A.2.4 Bayes' Rule

When two random variables are jointly distributed with the value of one known, the probability that the other takes a given value can be computed using *Bayes' rule*:

$$(A.18) \quad P(y|x) = \frac{P(x|y)P_Y(y)}{P_X(x)} = \frac{P(x|y)P_Y(y)}{\sum_y P(x|y)P_Y(y)}$$

Or, in words

$$(A.19) \quad \text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

Note that the denominator is obtained by summing (or integrating if  $y$  is continuous) the numerator over all possible  $y$  values. The “shape” of  $p(y|x)$  depends on the numerator with denominator as a normalizing factor to guarantee that  $p(y|x)$  sum to 1. Bayes' rule allows us to modify a prior probability into a posterior probability by taking information provided by  $x$  into account.

Bayes' rule inverts dependencies, allowing us to compute  $p(y|x)$  if  $p(x|y)$  is known. Suppose that  $y$  is the “cause” of  $x$ , like  $y$  going on summer vacation and  $x$  having a suntan. Then  $p(x|y)$  is the probability that someone who is known to have gone on summer vacation has a suntan. This is the *causal* (or predictive) way. Bayes' rule allows us a *diagnostic* approach by allowing us to compute  $p(y|x)$ : namely, the probability that someone who is known to have a suntan, has gone on summer vacation. Then  $p(y)$  is the general probability of anyone's going on summer vacation and  $p(x)$  is the probability that anyone has a suntan, including both those who have gone on summer vacation and those who have not.

### A.2.5 Expectation

*Expectation, expected value, or mean* of a random variable  $X$ , denoted by  $E[X]$ , is the average value of  $X$  in a large number of experiments:

$$(A.20) \quad E[X] = \begin{cases} \sum_i x_i P(x_i) & \text{if } X \text{ is discrete} \\ \int x p(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

It is a weighted average where each value is weighted by the probability that  $X$  takes that value. It has the following properties ( $a, b \in \mathfrak{R}$ ):

$$(A.21) \quad \begin{aligned} E[aX + b] &= aE[X] + b \\ E[X + Y] &= E[X] + E[Y] \end{aligned}$$

For any real-valued function  $g(\cdot)$ , the expected value is

$$(A.22) \quad E[g(X)] = \begin{cases} \sum_i g(x_i)P(x_i) & \text{if } X \text{ is discrete} \\ \int g(x)p(x)dx & \text{if } X \text{ is continuous} \end{cases}$$

A special  $g(x) = x^n$ , called the  $n$ th moment of  $X$ , is defined as

$$(A.23) \quad E[X^n] = \begin{cases} \sum_i x_i^n P(x_i) & \text{if } X \text{ is discrete} \\ \int x^n p(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

*Mean* is the first moment and is denoted by  $\mu$ .

### A.2.6 Variance

*Variance* measures how much  $X$  varies around the expected value. If  $\mu \equiv E[X]$ , the variance is defined as

$$(A.24) \quad \text{Var}(X) = E[(X - \mu)^2] = E[X^2] - \mu^2$$

Variance is the second moment minus the square of the first moment. Variance, denoted by  $\sigma^2$ , satisfies the following property ( $a, b \in \mathfrak{R}$ ):

$$(A.25) \quad \text{Var}(aX + b) = a^2 \text{Var}(X)$$

$\sqrt{\text{Var}(X)}$  is called the *standard deviation* and is denoted by  $\sigma$ . Standard deviation has the same unit as  $X$  and is easier to interpret than variance.

*Covariance* indicates the relationship between two random variables. If the occurrence of  $X$  makes  $Y$  more likely to occur, then the covariance is positive; it is negative if  $X$ 's occurrence makes  $Y$  less likely to happen and is 0 if there is no dependence.

$$(A.26) \quad \text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = E[XY] - \mu_X \mu_Y$$

where  $\mu_X \equiv E[X]$  and  $\mu_Y \equiv E[Y]$ . Some other properties are

$$\text{Cov}(X, Y) = \text{Cov}(Y, X)$$

$$\text{Cov}(X, X) = \text{Var}(X)$$

$$\text{Cov}(X + Z, Y) = \text{Cov}(X, Y) + \text{Cov}(Z, Y)$$

$$(A.27) \quad \text{Cov}\left(\sum_i X_i, Y\right) = \sum_i \text{Cov}(X_i, Y)$$

$$(A.28) \quad \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

$$(A.29) \quad \text{Var}\left(\sum_i X_i\right) = \sum_i \text{Var}(X_i) + \sum_i \sum_{j \neq i} \text{Cov}(X_i, X_j)$$

If  $X$  and  $Y$  are independent,  $E[XY] = E[X]E[Y] = \mu_X\mu_Y$  and  $\text{Cov}(X, Y) = 0$ . Thus if  $X_i$  are independent

$$(A.30) \quad \text{Var} \left( \sum_i X_i \right) = \sum_i \text{Var}(X_i)$$

*Correlation* is a normalized, dimensionless quantity that is always between  $-1$  and  $1$ :

$$(A.31) \quad \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

### A.2.7 Weak Law of Large Numbers

Let  $\mathcal{X} = \{X^t\}_{t=1}^N$  be a set of independent and identically distributed (iid) random variables each having mean  $\mu$  and a finite variance  $\sigma^2$ . Then for any  $\epsilon > 0$

$$(A.32) \quad P \left\{ \left| \frac{\sum_t X^t}{N} - \mu \right| > \epsilon \right\} \rightarrow 0 \text{ as } N \rightarrow \infty$$

That is, the average of  $N$  trials converges to the mean as  $N$  increases.

## A.3 Special Random Variables

There are certain types of random variables that occur so frequently that names are given to them.

### A.3.1 Bernoulli Distribution

A trial is performed whose outcome is either a “success” or a “failure.” The random variable  $X$  is a 0/1 indicator variable and takes the value 1 for a success outcome and is 0 otherwise.  $p$  is the probability that the result of trial is a success. Then

$$(A.33) \quad P\{X = 1\} = p \text{ and } P\{X = 0\} = 1 - p$$

which can equivalently be written as

$$(A.34) \quad P\{X = i\} = p^i(1 - p)^{1-i}, i = 0, 1$$

If  $X$  is Bernoulli, its expected value and variance are

$$(A.35) \quad E[X] = p, \text{ Var}(X) = p(1 - p)$$

### A.3.2 Binomial Distribution

If  $N$  identical independent Bernoulli trials are made, the random variable  $X$  that represents the number of successes that occurs in  $N$  trials is binomial distributed. The probability that there are  $i$  successes is

$$(A.36) \quad P\{X = i\} = \binom{N}{i} p^i (1-p)^{N-i}, i = 0 \dots N$$

If  $X$  is binomial, its expected value and variance are

$$(A.37) \quad E[X] = Np, \text{Var}(X) = Np(1-p)$$

### A.3.3 Multinomial Distribution

Consider a generalization of Bernoulli where instead of two states, the outcome of a random event is one of  $K$  mutually exclusive and exhaustive states, each of which has a probability of occurring  $p_i$  where  $\sum_{i=1}^K p_i = 1$ . Suppose that  $N$  such trials are made where outcome  $i$  occurred  $N_i$  times with  $\sum_{i=1}^K N_i = N$ . Then the joint distribution of  $N_1, N_2, \dots, N_K$  is multinomial:

$$(A.38) \quad P(N_1, N_2, \dots, N_K) = N! \prod_{i=1}^K \frac{p_i^{N_i}}{N_i!}$$

A special case is when  $N = 1$ ; only one trial is made. Then  $N_i$  are 0/1 indicator variables of which only one of them is 1 and all others are 0. Then equation A.38 reduces to

$$(A.39) \quad P(N_1, N_2, \dots, N_K) = \prod_{i=1}^K p_i^{N_i}$$

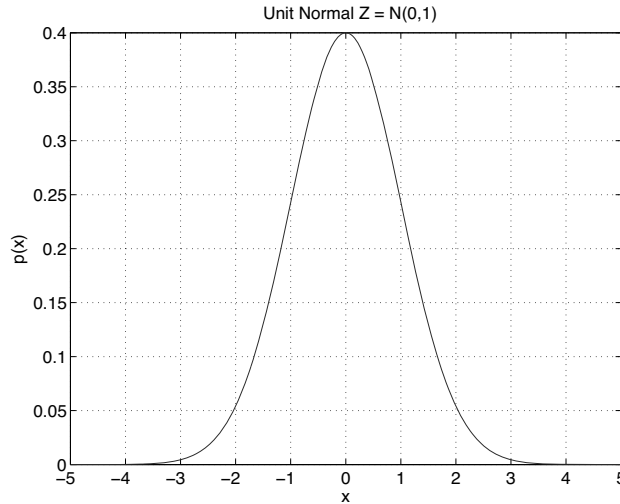
### A.3.4 Uniform Distribution

$X$  is uniformly distributed over the interval  $[a, b]$  if its density function is given by

$$(A.40) \quad p(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

If  $X$  is uniform, its expected value and variance are

$$(A.41) \quad E[X] = \frac{a+b}{2}, \text{Var}(X) = \frac{(b-a)^2}{12}$$



**Figure A.1** Probability density function of  $Z$ , the unit normal distribution.

### A.3.5 Normal (Gaussian) Distribution

$X$  is normal or Gaussian distributed with mean  $\mu$  and variance  $\sigma^2$ , denoted as  $\mathcal{N}(\mu, \sigma^2)$ , if its density function is

$$(A.42) \quad p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad -\infty < x < \infty$$

Many random phenomena obey the bell-shaped normal distribution, at least approximately, and many observations from nature can be seen as a continuous, slightly different versions of a typical value—that is probably why it is called the *normal* distribution. In such a case,  $\mu$  represents the typical value and  $\sigma$  defines how much instances vary around the prototypical value.

68.27 percent lie in  $(\mu - \sigma, \mu + \sigma)$ , 95.45 percent in  $(\mu - 2\sigma, \mu + 2\sigma)$ , and 99.73 percent in  $(\mu - 3\sigma, \mu + 3\sigma)$ . Thus  $P\{|x - \mu| < 3\sigma\} \approx 0.99$ . For practical purposes,  $p(x) \approx 0$  if  $x < \mu - 3\sigma$  or  $x > \mu + 3\sigma$ .  $Z$  is unit normal, namely,  $\mathcal{N}(0, 1)$  (see figure A.1), and its density is written as

$$(A.43) \quad p_Z(x) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{x^2}{2}\right]$$

If  $X \sim \mathcal{N}(\mu, \sigma^2)$  and  $Y = aX + b$ , then  $Y \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$ . The sum of independent normal variables is also normal with  $\mu = \sum_i \mu_i$  and  $\sigma^2 = \sum_i \sigma_i^2$ . If  $X$  is  $\mathcal{N}(\mu, \sigma^2)$ , then

$$(A.44) \quad \frac{X - \mu}{\sigma} \sim \mathcal{Z}$$

This is called z-normalization.

CENTRAL LIMIT  
THEOREM

Let  $X_1, X_2, \dots, X_N$  be a set of iid random variables all having mean  $\mu$  and variance  $\sigma^2$ . Then the *central limit theorem* states that for large  $N$ , the distribution of

$$(A.45) \quad X_1 + X_2 + \dots + X_N$$

is approximately  $\mathcal{N}(N\mu, N\sigma^2)$ . For example, if  $X$  is binomial with parameters  $(N, p)$ ,  $X$  can be written as the sum of  $N$  Bernoulli trials and  $(X - Np)/\sqrt{Np(1-p)}$  is approximately unit normal.

Central limit theorem is also used to generate normally distributed random variables on computers. Programming languages have subroutines that return uniformly distributed (pseudo-)random numbers in the range  $[0, 1]$ . When  $U_i$  are such random variables,  $\sum_{i=1}^{12} U_i - 6$  is approximately  $\mathcal{Z}$ .

Let us say  $X^t \sim \mathcal{N}(\mu, \sigma^2)$ . The estimated sample mean

$$(A.46) \quad m = \frac{\sum_{t=1}^N X^t}{N}$$

is also normal with mean  $\mu$  and variance  $\sigma^2/N$ .

### A.3.6 Chi-Square Distribution

If  $Z_i$  are independent unit normal random variables, then

$$(A.47) \quad X = Z_1^2 + Z_2^2 + \dots + Z_n^2$$

is chi-square with  $n$  degrees of freedom, namely,  $X \sim \mathcal{X}_n^2$ , with

$$(A.48) \quad E[X] = n, \text{ Var}(X) = 2n$$

When  $X^t \sim \mathcal{N}(\mu, \sigma^2)$ , the estimated sample variance is

$$(A.49) \quad S^2 = \frac{\sum_t (X^t - m)^2}{N - 1}$$

and we have

$$(A.50) \quad (N - 1) \frac{S^2}{\sigma^2} \sim \mathcal{X}_{N-1}^2$$

It is also known that  $m$  and  $S^2$  are independent.

**A.3.7  $t$  Distribution**

If  $Z \sim \mathcal{Z}$  and  $X \sim \mathcal{X}_n^2$  are independent, then

$$(A.51) \quad T_n = \frac{Z}{\sqrt{X/n}}$$

is  $t$ -distributed with  $n$  degrees of freedom with

$$(A.52) \quad E[T_n] = 0, n > 1, \text{Var}(T_n) = \frac{n}{n-2}, n > 2$$

Like the unit normal density,  $t$  is symmetric around 0. As  $n$  becomes larger,  $t$  density becomes more and more like the unit normal, the difference being that  $t$  has thicker tails, indicating greater variability than does normal.

**A.3.8  $F$  Distribution**

If  $X_1 \sim \mathcal{X}_n^2$  and  $X_2 \sim \mathcal{X}_m^2$  are independent chi-square random variables with  $n$  and  $m$  degrees of freedom, respectively,

$$(A.53) \quad F_{n,m} = \frac{X_1/n}{X_2/m}$$

is  $F$ -distributed with  $n$  and  $m$  degrees of freedom with

$$(A.54) \quad E[F_{n,m}] = \frac{m}{m-2}, m > 2, \text{Var}(F_{n,m}) = \frac{m^2(2m+2n-4)}{n(m-2)^2(m-4)}, m > 4$$

**A.4 References**

- Casella, G., and R. L. Berger. 1990. *Statistical Inference*. Belmont, CA: Duxbury.
- Ross, S. M. 1987. *Introduction to Probability and Statistics for Engineers and Scientists*. New York: Wiley.



## *Index*

- 0/1 loss function, 51
- 5×2
  - cross-validation, 488
  - cv paired  $F$  test, 503
  - cv paired  $t$  test, 503
- Active learning, 360
- AdaBoost, 431
- Adaptive resonance theory, 285
- Additive models, 180
- Agglomerative clustering, 157
- AIC, *see* Akaike's information criterion
- Akaike's information criterion, 81
- Alignment, 324
- Analysis of variance, 504
- Anchor, 291
- ANOVA, *see* Analysis of variance
- Approximate normal test, 500
- Apriori algorithm, 56
- Area under the curve, 491
- ART, *see* Adaptive resonance theory
- Artificial neural networks, 233
- Association rule, 4, 55
- Attribute, 87
- AUC, *see* Area under the curve
- Autoassociator, 268
  
- Backpropagation, 250
  - through time, 272
- Backup, 456
  
- Backward selection, 111
- Backward variable, 372
- Bag of words, 102, 324
- Bagging, 430
- Base-learner, 419
- Basis function, 211
  - cooperative vs. competitive, 297
  - for a kernel, 352
  - normalization, 295
- Basket analysis, 55
- Batch learning, 251
- Baum-Welch algorithm, 376
- Bayes' ball, 402
- Bayes' classifier, 51
- Bayes' estimator, 68
- Bayes' rule, 49, 521
- Bayesian information criterion, 81
- Bayesian model combination, 426
- Bayesian model selection, 82
- Bayesian networks, 387
- Belief networks, 387
- Belief state, 465
- Bellman's equation, 452
- Beta distribution, 345
- Between-class scatter matrix, 130
- Bias, 65
- Bias unit, 237
- Bias/variance dilemma, 78
- BIC, *see* Bayesian information criterion
- Binary split, 187

- Binding, 202
- Binomial test, 499
- Biometrics, 441
- Blocking, 482
- Bonferroni correction, 508
- Boosting, 431
- Bootstrap, 489
  
- C4.5, 191
- C4.5Rules, 197
- CART, 191, 203
- Cascade correlation, 264
- Cascading, 438
- Case-based reasoning, 180
- Causality, 396
  - causal graph, 388
- Central limit theorem, 526
- Class
  - confusion matrix, 493
  - likelihood, 50
- Classification, 5
  - likelihood- vs.
    - discriminant-based, 209
- Classification tree, 188
- Clique, 411
- Cluster, 144
- Clustering, 11
  - agglomerative, 157
  - divisive, 157
  - hierarchical, 157
  - online, 281
- Code word, 146
- Codebook vector, 146
- Coefficient of determination (of regression), 76
- Color quantization, 145
- Common principal components, 119
- Competitive basis functions, 297
- Competitive learning, 280
- Complete-link clustering, 158
- Component density, 144
  
- Compression, 8, 146
- Condensed nearest neighbor, 173
- Conditional independence, 389
- Confidence interval
  - one-sided, 495
  - two-sided, 494
- Confidence of an association rule, 55
- Conjugate prior, 344
- Connection weight, 237
- Contingency table, 501
- Correlation, 89
- Cost-sensitive learning, 478
- Coupled HMM, 400
- Covariance function, 356
- Covariance matrix, 88
- Credit assignment, 448
- Critic, 448
- CRM, *see* Customer relationship management
- Cross-entropy, 221
- Cross-validation, 40, 80, 486
  - $5 \times 2$ , 488
  - $K$ -fold, 487
- Curse of dimensionality, 170
- Customer relationship management, 155
- Customer segmentation, 155
  
- $d$ -separation, 402
- Decision node, 185
- Decision region, 53
- Decision tree, 185
  - multivariate, 202
  - omnivariate, 205
  - soft, 305
  - univariate, 187
- Delve repository, 17
- Dendrogram, 158
- Density estimation, 11
- Dichotomizer, 53
- Diffusion kernel, 325

- Dimensionality reduction
  - nonlinear, 269
- Directed acyclic graph, 387
- Dirichlet distribution, 344
- Discount rate, 451
- Discriminant, 5
  - function, 53
  - linear, 97
  - quadratic, 95
- Discriminant adaptive nearest neighbor, 172
- Discriminant-based classification, 209
- Distributed vs. local representation, 156, 289
- Diversity, 420
- Divisive clustering, 157
- Document categorization, 102
- Doubt, 26
- Dual representation, 337, 352
- Dynamic classifier selection, 435
- Dynamic graphical models, 415
- Dynamic node creation, 264
- Dynamic programming, 453
  
- Early stopping, 223, 258
- ECOC, 327, *see* Error-correcting output codes
- Edit distance, 324
- Eigendigits, 118
- Eigenfaces, 118
- Eligibility trace, 459
- EM, *see* Expectation-Maximization
- Emission probability, 367
- Empirical error, 24
- Empirical kernel map, 324
- Ensemble, 424
- Ensemble selection, 437
- Entropy, 188
- Episode, 451
- Epoch, 251
- Error
  - type I, 497
  - type II, 497
- Error-correcting output codes, 427
- Euclidean distance, 98
- Evidence, 50
- Example, 87
- Expectation-Maximization, 150
  - supervised, 299
- Expected error, 476
- Expected utility, 54
- Experiment
  - design, 478
  - factorial, 481
  - strategies, 480
- Explaining away, 393
- Extrapolation, 35
  
- FA, *see* Factor analysis
- Factor analysis, 120
- Factor graph, 412
- Factorial HMM, 400
- Feature, 87
  - extraction, 110
  - selection, 110
- Finite-horizon, 451
- First-order rule, 201
- Fisher kernel, 325
- Fisher's linear discriminant, 129
- Flexible discriminant analysis, 120
- Floating search, 112
- Foil, 199
- Forward selection, 110
- Forward variable, 370
- Forward-backward procedure, 370
- Fuzzy  $k$ -means, 160
- Fuzzy membership function, 295
- Fuzzy rule, 295
  
- Gamma distribution, 347
- Gamma function, 344
- Gaussian prior, 349
- Generalization, 24, 39

- Generalized linear models, 230
- Generative model, 342, 397
- Generative topographic mapping, 306
- Geodesic distance, 133
- Gini index, 189
- Gradient descent, 219
  - stochastic, 241
- Gradient vector, 219
- Gram matrix, 321
- Graphical models, 387
- Group, 144
- GTM, *see* Generative topographic mapping
- Hamming distance, 171
- Hebbian learning, 283
- Hidden layer, 246
- Hidden Markov model, 367, 398
  - coupled, 400
  - factorial, 400
  - input-output, 379, 400
  - left-to-right, 380
  - switching, 400
- Hidden variables, 57, 396
- Hierarchical clustering, 157
- Hierarchical cone, 260
- Hierarchical mixture of experts, 304
- Higher-order term, 211
- Hinge loss, 317
- Hint, 261
- Histogram, 165
- HMM, *see* Hidden Markov model
- Hybrid learning, 291
- Hypothesis, 23
  - class, 23
  - most general, 24
  - most specific, 24
- Hypothesis testing, 496
- ID3, 191
- IF-THEN rules, 197
- lid (independent and identically distributed), 41
- Ill-posed problem, 38
- Impurity measure, 188
- Imputation, 89
- Independence, 388
- Inductive bias, 38
- Inductive logic programming, 202
- Infinite-horizon, 451
- Influence diagrams, 414
- Information retrieval, 491
- Initial probability, 364
- Input, 87
- Input representation, 21
- Input-output HMM, 379, 399
- Instance, 87
- Instance-based learning, 164
- Interest of an association rule, 55
- Interpolation, 35
- Interpretability, 197
- Interval estimation, 493
- Irep, 199
- Isometric feature mapping, 133
- Job shop scheduling, 471
- Junction tree, 410
- $K$ -armed bandit, 449
- $K$ -fold
  - cross-validation, 487
  - cv paired  $t$  test, 502
- $k$ -means clustering, 147
  - fuzzy, 160
  - online, 281
- $k$ -nearest neighbor
  - classifier, 172
  - density estimate, 169
  - smoother, 177
- $k$ -nn, *see*  $k$ -nearest neighbor
- Kalman filter, 400
- Karhunen-Loève expansion, 119
- Kernel estimator, 167

- Kernel function, 167, 320, 353
- Kernel PCA, 336
- Kernel smoother, 176
- kernelization, 321
- Knowledge extraction, 8, 198, 295
- Kolmogorov complexity, 82
- Kruskal-Wallis test, 511
  
- Laplace approximation, 354
- Laplacian prior, 350
- lasso, 352
- Latent factors, 120
- Lateral inhibition, 282
- LDA, *see* Linear discriminant analysis
- Leader cluster algorithm, 148
- Leaf node, 186
- Learning automata, 471
- Learning vector quantization, 300
- Least square difference test, 507
- Least squares estimate, 74
- Leave-one-out, 487
- Left-to-right HMM, 380
- Level of significance, 497
- Levels of analysis, 234
- Lift of an association rule, 55
- Likelihood, 62
- Likelihood ratio, 58
- Likelihood-based classification, 209
- Linear classifier, 97, 216
- Linear discriminant, 97, 210
- Linear discriminant analysis, 128
- Linear dynamical system, 400
- Linear opinion pool, 424
- Linear regression, 74
  - multivariate, 103
- Linear separability, 215
- Local representation, 288
- Locally linear embedding, 135
- Locally weighted running line smoother, 177
  
- Loess, *see* Locally weighted running line smoother
- Log likelihood, 62
- Log odds, 58, 218
- Logistic discrimination, 220
- Logistic function, 218
- Logit, 218
- Loss function, 51
- LSD, *see* Least square difference test
- LVQ, *see* Learning vector quantization
  
- Mahalanobis distance, 90
- Margin, 25, 311, 433
- Markov decision process, 451
- Markov mixture of experts, 379
- Markov model, 364
  - hidden, 367
  - learning, 366, 375
  - observable, 365
- Markov random field, 410
- Max-product algorithm, 413
- Maximum a posteriori (MAP)
  - estimate, 68, 343
- Maximum likelihood estimation, 62
- McNemar's test, 501
- MDP, *see* Markov decision process
- MDS, *see* Multidimensional scaling
- Mean square error, 65
- Mean vector, 88
- Memory-based learning, 164
- Minimum description length, 82
- Mixture components, 144
- Mixture density, 144
- Mixture of experts, 301, 434
  - competitive, 304
  - cooperative, 303
  - hierarchical, 305
  - Markov, 379, 400
- Mixture of factor analyzers, 155
- Mixture of mixtures, 156

- Mixture of probabilistic principal component analyzers, 155
- Mixture proportion, 144
- MLE, *see* Maximum likelihood estimation
- Model combination
  - multiexpert, 423
  - multistage, 423
- Model selection, 38
- MoE, *see* Mixture of experts
- Momentum, 257
- Moralization, 411
- Multidimensional scaling, 125
  - nonlinear, 287
  - using MLP, 269
- Multilayer perceptrons, 246
- Multiple comparisons, 507
- Multiple kernel learning, 326, 442
- Multivariate linear regression, 103
- Multivariate polynomial regression, 104
- Multivariate tree, 202
- Naive Bayes' classifier, 397
  - discrete inputs, 102
  - numeric inputs, 97
- Naive estimator, 166
- Nearest mean classifier, 98
- Nearest neighbor classifier, 172
  - condensed, 173
- Negative examples, 21
- Neuron, 233
- No Free Lunch Theorem, 477
- Noise, 30
- Noisy OR, 409
- Nonparametric estimation, 163
- Nonparametric tests, 508
- Null hypothesis, 497
- Observable Markov model, 365
- Observable variable, 48
- Observation, 87
- Observation probability, 367
- OC1, 203
- Occam's razor, 32
- Off-policy, 458
- Omnivariate decision tree, 205
- On-policy, 458
- One-class classification, 333
- One-sided confidence interval, 495
- One-sided test, 498
- Online  $k$ -means, 281
- Online learning, 241
- Optimal policy, 452
- Optimal separating hyperplane, 311
- Outlier detection, 9, 333
- Overfitting, 39, 79
- Overtraining, 258
- PAC, *see* Probably approximately correct
- Paired test, 501
- Pairing, 482
- Pairwise separation, 216, 428
- Parallel processing, 236
- Partially observable Markov decision process, 464
- Parzen windows, 167
- Pattern recognition, 6
- PCA, *see* Principal components analysis
- Pedigree, 400
- Perceptron, 237
- Phone, 381
- Phylogenetic tree, 398
- Piecewise approximation
  - constant, 248, 300
  - linear, 301
- Policy, 451
- Polychotomizer, 53
- Polynomial regression, 75
  - multivariate, 104
- Polytree, 407
- POMDP, *see* Partially observable Markov decision process

- Positive examples, 21
- Posterior probability distribution, 341
- Posterior probability of a class, 50
- Posterior probability of a parameter, 67
- Posthoc testing, 507
- Postpruning, 194
- Potential function, 212, 411
- Power function, 498
- Precision
  - in information retrieval, 492
  - reciprocal of variance, 347
- Predicate, 201
- Prediction, 5
- Prepruning, 194
- Principal components analysis, 113
- Prior knowledge, 294
- Prior probability distribution, 341
- Prior probability of a class, 50
- Prior probability of a parameter, 67
- Probabilistic networks, 387
- Probabilistic PCA, 123
- Probably approximately correct learning, 29
- Probit function, 355
- Product term, 211
- Projection pursuit, 274
- Proportion of variance, 116
- Propositional rule, 201
- Pruning
  - postpruning, 194
  - prepruning, 194
  - set, 194
- Q learning, 458
- Quadratic discriminant, 95, 211
- Quantization, 146
- Radial basis function, 290
- Random Subspace, 421
- Randomization, 482
- RBF, *see* Radial basis function
- Real time recurrent learning, 272
- Recall, 492
- Receiver operating characteristics, 490
- Receptive field, 288
- Reconstruction error, 119, 146
- Recurrent network, 271
- Reference vector, 146
- Regression, 9, 35
  - linear, 74
  - polynomial, 75
  - polynomial multivariate, 104
  - robust, 329
- Regression tree, 192
- Regressogram, 175
- Regularization, 80, 266
- Regularized discriminant analysis, 100
- Reinforcement learning, 13
- Reject, 34, 52
- Relative square error, 76
- Replication, 482
- Representation, 21
  - distributed vs. local, 288
- Response surface design, 481
- Ridge regression, 266, 350
- Ripper, 199
- Risk function, 51
- Robust regression, 329
- ROC, *see* Receiver operating characteristics
- RSE, *see* Relative square error
- Rule
  - extraction, 295
  - induction, 198
  - pruning, 198
- Rule support, 198
- Rule value metric, 199
- Running smoother
  - line, 177
  - mean, 175

- Sammon mapping, 128
  - using MLP, 269
- Sammon stress, 128
- Sample, 48
  - correlation, 89
  - covariance, 89
  - mean, 89
- Sarsa, 458
  - Sarsa( $\lambda$ ), 461
- Scatter, 129
- Scree graph, 116
- Self-organizing map, 286
- Semiparametric density estimation, 144
- Sensitivity, 493
- Sensor fusion, 421
- Sequential covering, 199
- Sigmoid, 218
- Sign test, 509
- Single-link clustering, 157
- Slack variable, 315
- Smoother, 174
- Smoothing splines, 178
- Soft count, 376
- Soft error, 315
- Soft weight sharing, 267
- Softmax, 224
- SOM, *see* Self-organizing map
- Spam filtering, 103
- Specificity, 493
- Spectral decomposition, 115
- Speech recognition, 380
- Sphere node, 203
- Stability-plasticity dilemma, 281
- Stacked generalization, 435
- Statlib repository, 17
- Stochastic automaton, 364
- Stochastic gradient descent, 241
- Stratification, 487
- Strong learner, 431
- Structural adaptation, 263
- Structural risk minimization, 82
- Subset selection, 110
- Sum-product algorithm, 412
- Supervised learning, 9
- Support of an association rule, 55
- Support vector machine, 313
- SVM, *see* Support vector machine
- Switching HMM, 400
- Synapse, 234
- Synaptic weight, 237
- $t$  distribution, 495
- $t$  test, 498
- Tangent prop, 263
- TD, *see* Temporal difference
- Template matching, 98
- Temporal difference, 455
  - learning, 458
  - TD(0), 459
  - TD-Gammon, 471
- Test set, 40
- Threshold, 212
  - function, 238
- Time delay neural network, 270
- Topographical map, 287
- Transition probability, 364
- Traveling salesman problem, 306
- Triple trade-off, 39
- Tukey's test, 512
- Two-sided confidence interval, 494
- Two-sided test, 497
- Type 2 maximum likelihood
  - procedure, 360
- Type I error, 497
- Type II error, 497
- UCI repository, 17
- Unbiased estimator, 65
- Underfitting, 39, 79
- Unfolding in time, 272
- Unit normal distribution, 493
- Univariate tree, 187
- Universal approximation, 248

Unobservable variable, 48  
Unstable algorithm, 430  
Utility function, 54  
Utility theory, 54

Validation set, 40  
Value iteration, 453  
Value of information, 464, 469  
Vapnik-Chervonenkis (VC)  
  dimension, 27  
Variance, 66  
Vector quantization, 146  
  supervised, 300  
Version space, 24  
Vigilance, 285  
Virtual example, 262  
Viterbi algorithm, 374  
Voronoi tessellation, 172  
Voting, 424

Weak learner, 431  
Weight  
  decay, 263  
  sharing, 260  
  sharing soft, 267  
  vector, 212  
Wilcoxon signed rank test, 511  
Winner-take-all, 280  
Within-class scatter matrix, 130  
Wrappers, 138

*z*, *see* Unit normal distribution  
*z*-normalization, 91, 526  
Zero-one loss, 51

## **Adaptive Computation and Machine Learning**

Thomas Dietterich, Editor  
Christopher Bishop, David Heckerman, Michael Jordan, and Michael  
Kearns, Associate Editors

*Bioinformatics: The Machine Learning Approach*, Pierre Baldi and Søren  
Brunak

*Reinforcement Learning: An Introduction*, Richard S. Sutton and Andrew  
G. Barto

*Graphical Models for Machine Learning and Digital Communication*,  
Brendan J. Frey

*Learning in Graphical Models*, Michael I. Jordan

*Causation, Prediction, and Search*, second edition, Peter Spirtes, Clark  
Glymour, and Richard Scheines

*Principles of Data Mining*, David Hand, Heikki Mannila, and Padhraic  
Smyth

*Bioinformatics: The Machine Learning Approach*, second edition, Pierre  
Baldi and Søren Brunak

*Learning Kernel Classifiers: Theory and Algorithms*, Ralf Herbrich

*Learning with Kernels: Support Vector Machines, Regularization,  
Optimization, and Beyond*, Bernhard Schölkopf and Alexander J. Smola

*Introduction to Machine Learning*, Ethem Alpaydm

*Gaussian Processes for Machine Learning*, Carl Edward Rasmussen and  
Christopher K. I. Williams

*Semi-Supervised Learning*, Olivier Chapelle, Bernhard Schölkopf, and  
Alexander Zien, Eds.

*The Minimum Description Length Principle*, Peter D. Grünwald

*Introduction to Statistical Relational Learning*, Lise Getoor and Ben Taskar, Eds.

*Probabilistic Graphical Models: Principles and Techniques*, Daphne Koller and Nir Friedman

*Introduction to Machine Learning*, second edition, Ethem Alpaydm