Thus, the tree can be automatically constructed by searching, for each node, the question that renders the maximum entropy decrease. Alternatively, complex questions can be formed for each node for improved splitting.

When we grow the tree, it needs to be pruned using cross-validation as discussed in Chapter 4. When the algorithm terminates, the leaf nodes of the tree represent the senones to be used. Figure 9.6 shows an example tree we built to classify the second state of all |k|triphones seen in a training corpus. After the tree is built, it can be applied to the second state of any /k/ triphone, thanks to the generalizability of the binary tree and the general linguistic questions. Figure 9.6 indicates that the second state of the /k/ triphone in welcome is mapped to the second senone, no matter whether this triphone occurs in the training corpus or not.



Figure 9.5 State-based vs. model-based clustering. These two models are very similar, as both the first and the second output distributions are almost identical. The key difference is the output distribution of the third state. If we measure the overall model similarity, which is often based on the accumulative output distribution similarities of all states, these two models may be clustered, leading to a very inaccurate distribution for the last state. If we cluster output distributions at state level, we can cluster the first two output distributions while leaving the last ones intact, leading to more accurate representations.

## Phonetic Modeling—Selecting Appropriate Units

Questions	Phones in Each Question Category
Aspseg	hh
Sil	sil
Alvstp	d t
Dental	dh th
Labstp	b p
Liquid	lr
Lw	lw
S/Sh	s sh
Sylbic	er axr
Velstp	gk
Affric	ch jh
Lqgl-B	lrw
Nasal	m n ng
Retro	r er axr
Schwa	ax ix axr
Velar	ng g k
Fric2	th s sh f
Fric3	dh z zh v
Lągl	lrwy
S/Z/Sh/Zh	s z sh zh
Wglide	uw aw ow w
Labial	w m b p v
Palatl	y ch jh sh zh
Yglide	iy ay ey oy y
High	ih ix iy uh uw y
Lax	eh ih ix uh ah ax
Low	ae aa ao aw ay oy
Orstp2	p t k
Orstp3	b d g
Alvelr	ndtsz
Diph	uw aw ay ey iy ow oy
Fricl	dh th s sh z zh v f
Round	uh ao uw ow oy w axr er
Frnt-R	ae eh ih ix iy ey ah ax y aw
Tense	iy ey ae uw ow aa ao ay oy aw
Back-L	uh ao uw ow aa er axr l r w aw
Frnt-L	ae eh ih ix iy ey ah ax y oy ay
Back-R	uh ao uw ow aa er axr oy l r w ay
Orstp 1	bdgptkchjh
Vowel	ae eh ih ix iy uh ah ax aa ao uw uw uy ey on oy or an uh ay uy ay ar aw lr w y
Son	ae eh ih ix iy ey ah ax oy ay un ao uw ow au or au w v er axr m
Voiced	ae eh ih ix iy uh ah ax aa ao uw aw ay ey on oy thay a
	n ng ih b d dh g v z zh

Table 9.3 Some example questions used in building senone trees.



Figure 9.6 A decision tree for classifying the second state of K-triphone HMMs [48].

In practice, senone models significantly reduce the word recognition error rate in comparison with model-based clustered triphone models, as illustrated in Table 9.4. It is the senonic model's significant reduction of the overall system parameters that enables the continuous mixture HMMs to perform well for large-vocabulary speech recognition [56].

**Table 9.4** Relative error reductions for different modeling units. The reduction is relative to that of the preceding row.

Units	Relative Error Reductions
Context-independent phone	Baseline
Context-dependent phone	+25%
Clustered triphone	+15%
Senone	+24%

#### 9.4.4. Lexical Baseforms

When appropriate subword units are used, we must have the correct pronunciation for each word so that concatenation of subword units can accurately represent the word to be recognized. The dictionary represents the standard pronunciation used as a starting point for building a workable speech recognition system. We also need to provide alternative pronunciations for words such as *tomato* that may have very different pronunciations. For example, the COMLEX dictionary from LDC has about 90,000 baseforms that cover most words used in many years of *The Wall Street Journal*. The CMU Pronunciation Dictionary, which was optimized for continuous speech recognition, has about 100,000 baseforms.

Amazon/VB Assets Exhibit 1012 Page 462

#### Phonetic Modeling-Selecting Appropriate Units

In continuous speech recognition, we must also use phonologic rules to modify interword pronunciations or to have reduced sounds. Assimilation is a typical coarticulation phenomenon—a change in a segment to make it more like a neighboring segment. Typical examples include phrases such as *did you /d ih jh y ah/, set you /s eh ch er/, last year / l ae s ch iy r/, because you've /b iy k ah zh uw v/,* etc. Deletion is also common in continuous speech. For example, /t/ and /d/ are often deleted before a consonant. Thus, in conversational speech, you may find examples like *find him /f ay n ix m/, around this /ix r aw n ih s/,* and *Let me in /l eh m eh n/.* 

Dictionaries often don't include proper names. For example, the 20,000 names included in the COMLEX dictionary are a small fraction of 1–2 million names in the USA. To deal with these new words, we often have to derive their pronunciation automatically. These new words have to be added on the fly, either by the user or through an interface from speech-aware applications. Unlike Spanish or Italian, rule-based letter-to-sound (LTS) conversion for English is often impractical, since so many words in English don't follow phonological rules. A trainable LTS converter is attractive, since its performance can be improved by constantly learning from examples so that it can generalize rules for the specific task. Trainable LTS converters can be based on neural networks, HMMs, or the CART described in Chapter 4. In practice, CART-based LTS has a very accurate performance [10, 61, 71, 89].

When CART is used, the basic YES-NO question for LTS conversion looks like: Is the second right letter 'p'? or: Is the first left output phone /ay/? The question for letters and phones can be on either the left or the right side. The range of question positions should be long enough to cover the most important phonological variations. Empirically, a 10-letter window (5 for left letter context and 5 for right letter context) and 3-phone window context is generally sufficient. A primitive set of questions can include all the singleton questions about each letter or phone identity. If we allow the node to have a complex question—that is, a combination of primitive questions—the depth of the tree can be greatly reduced and performance improved. For example, a complex question: Is the second left letter 't' and the first left letter 'n'? can capture o in the common suffix tion and convert it to the correct phone. Complex questions can also alleviate possible data-fragmentation problems caused by the greedy nature of the CART algorithm.

Categorical questions can be formed in both the letter and phone domains with our common linguistic knowledge. For example, the most often used set includes the letter or phone clusters for vowels, consonants, nasals, liquids, fricatives, and so on. In growing the decision tree, the context distance also plays a major role in the overall quality. It is very important to weight the entropy reduction according to the distance (either letter or phoneme) to avoid overgeneralization, which forces the tree to look more carefully at the *nearby* context than at the *far-away* context. Each leaf of the tree has a probability distribution for letter-to-phoneme mapping.

There are a number of ways to improve the effectiveness of the decision tree. First, pruning controls the tree's depth. For example, certain criteria have to be met for a node to be split. Typically splitting requires a minimum number of counts and a minimum entropy reduction. Second, the distribution at the leaves can be smoothed. For example, a leaf distribution

# Acoustic Modeling

bution can be interpolated with the distributions of its ancestor nodes using deleted. interpolation. Finally, we can partition the training data and build multiple trees with different prediction capabilities. These trees accommodate different phonological rules with different language origins.

When the decision tree is used to derive the phonetic pronunciation, the phonetic conversion error is about 8% for the Wall Street Journal newspaper text corpora [61]. These errors can be broadly classified into two categories. The first includes errors of proper nouns and foreign words. For example, Pacino can be mistakenly converted to  $/p \ ax \ s \ iy \ n \ ow'$  instead of  $/p \ ax \ ch \ iy \ n \ ow'$ . The second category includes generalization errors. For example, shier may be converted to  $/sh \ ih \ r'$  instead of the correct pronunciation  $/sh \ ay \ r'$  if the word cashier  $/k \ ae \ sh \ ih \ r'$  appears in the training data. The top three phone confusion pairs are /ix/ax/, /dx/t/, and /ae/ax/. The most confusing pair is /ix/ax/. This is not surprising, because /ix/ax/ is among the most inconsistent transcriptions in most of the published dictionaries. There is no consensus for /ix/ax/ transcription among phoneticians.

Although automatic LTS conversion has a reasonable accuracy, it is hardly practical if you don't use an exception dictionary. This is especially true for proper nouns. In practice, you can often ask the person who knows how to pronounce the word to either speak or write down the correct phonetic pronunciation, updating the exception dictionary if the correct one disagrees with what the LTS generates. When acoustic examples are available, you can use the decision tree to generate multiple results and use these results as a language model to perform phone recognition on the acoustic examples. The best overall acoustic and LTS probability can be used as the most likely candidate in the exception dictionary. Since there may be many ways to pronounce a word, you can keep multiple pronunciations in the dictionary with a probability for each possible one. If the pronunciation probability is inaccurate, an increase in multiple pronunciations essentially increases the size and confusion of the vocabulary, leading to increased speech recognition error rate.

Even if you have accurate phonetic baseforms, pronunciations in spontaneous speech differ significantly from the standard baseform. Analysis of manual phonetic transcription of conversational speech reveals a large number (> 20%) of cases of genuine ambiguity: instances where human labelers disagree on the identity of the surface form [95]. For example, the word because has more than 15 different pronunciation variations, such as */b ix k ah z/*, */k ah z/*, */k ax z/*, */b ix k ax z/*, */b ax k ah z/*, */b ih k ah z/*, */k x/k ix z/*, */k ih z//b iy k ah z/*, */b iy k ah z/*, */a x z/*, *k in z/*, *k ax z/*, *k in z/*, *k* 

To incorporate widespread pronunciations, we can use a probabilistic finite state machine to model each word's pronunciation variations, as shown in Figure 9.7. The probability with each arc indicates how likely that path is to be taken, with all the arcs that leave a node summing to 1. As with HMMs, these weights can be estimated from a real corpus for

# Acoustic Modeling-Scoring Acoustic Features

improved speech recognition [20, 85, 102, 103, 110]. In practice, the relative error reduction of using probabilistic finite state machines is very modest (5-10%).



Figure 9.7 A possible pronunciation network for word *tomato*. The vowel leyl is more likely to flap, thereby having a higher transition probability into ldx/.

### 9.5. ACOUSTIC MODELING—SCORING ACOUSTIC FEATURES

After feature extraction, we have a sequence of feature vectors,  $\mathbf{X}$ , such as the MFCC vector, as our input data. We need to estimate the probability of these acoustic features, given the word or phonetic model,  $\mathbf{W}$ , so that we can recognize the input data for the correct word. This probability is referred to as acoustic probability,  $P(\mathbf{X} \mid \mathbf{W})$ . In this section we focus our discussion on the HMM. As discussed in Chapter 8, it is the most successful method for acoustic modeling. Other emerging techniques are discussed in Section 9.8.

#### 9.5.1. Choice of HMM Output Distributions

As discussed in Chapter 8, you can use discrete, continuous, or semicontinuous HMMs. When the amount of training data is sufficient, parameter tying becomes unnecessary. A continuous model with a large number of mixtures offers the best recognition accuracy, although its computational complexity also increases linearly with the number of mixtures. On the other hand, the discrete model is computationally efficient, but has the worst performance among the three models. The semicontinuous model provides a viable alternative between system robustness and trainability.

When either the discrete or the semicontinuous HMM is employed, it is helpful to use multiple codebooks for a number of features for significantly improved performance. Each codebook then represents a set of different speech parameters. One way to combine these multiple output observations is to assume that they are independent, computing the output probability as the product of the output probabilities of each codebook. For example, the semicontinuous HMM output probability of multiple codebooks can be computed as the product of each codebook:

Acoustic Modeling

$$b_i(\mathbf{x}) = \prod_m \sum_{k=1}^{L^m} f^m(\mathbf{x}^m \mid o_k^m) b_i^m(o_k^m)$$
(9.12)

where superscript m denotes the codebook-m related parameters. Each codebook consists of  $L^{m}$ -mixture continuous density functions.

Following our discussion in Chapter 8, the re-estimation algorithm for the multiplecodebook-based HMM could be extended. Since multiplication of the output probability density of each codebook leads to several independent terms in the Q-function, for codebook  $m, \zeta_i(j,k^m)$  can be modified as follows:

$$\zeta_{I}(j,k^{m}) = \frac{\sum_{i} \alpha_{I-1}(i) a_{ij} b_{j}^{m}(k^{m}) f^{m}(\mathbf{x}_{I} \mid v_{k}^{m}) \prod_{m \neq n-k} \sum_{k} b_{j}^{n}(k^{n}) f^{n}(\mathbf{x}_{I} \mid v_{k}^{n}) \beta_{I}(j)}{\sum_{k} \alpha_{T}^{m}(k)}$$
(9.13)

Other intermediate probabilities can also be computed in a manner similar to what we discussed in Chapter 8.

Multiple codebooks can dramatically increase the representation power of the VQ codebook and can substantially improve speech recognition accuracy. You can typically build a codebook for  $\mathbf{c}_k$ ,  $\Delta \mathbf{c}_k$ , and  $\Delta \Delta \mathbf{c}_k$ , respectively. As energy has a very different dynamic range, you can further improve the performance by building a separate codebook for  $\mathbf{c}_k[0]$ ,  $\Delta \mathbf{c}_k[0]$ , and  $\Delta \Delta \mathbf{c}_k[0]$ . In comparison to building a single codebook for  $\mathbf{x}_k$  as illustrated in Eq. (9.6), the multiple-codebook system can reduce the error rate by more than 10%.

In practice, the most important parameter for the output probability distribution is the number of mixtures or the size of the codebooks. When there are sufficient training data, relative error reductions with respect to the discrete HMM are those shown in Figure 9.8.



Figure 9.8 Continuous speaker-independent word recognition error rates of the discrete HMM (DHMM), SCHMM, and the continuous HMM (CHMM) with respect to the training set sizes (thousands of training sentences). Both the DHMM and SCHMM have multiple codebooks. The CHMM has 20 mixture diagonal Gaussian density functions.

Amazon/VB Assets Exhibit 1012 Page 466

### Acoustic Modeling-Scoring Acoustic Features

As you can see from Figure 9.8, the SCHMM offers improved accuracy in comparison with the discrete HMM or the continuous HMM when the amount of training data is limited. When we increase the training data size, the continuous mixture density HMM starts to outperform both the discrete and the semicontinuous HMM, since the need to share model parameters becomes less critical.

Performance is also a function of the number of mixtures. With a small number of mixtures, the continuous HMM lacks the modeling power and it actually performs worse than the discrete HMM across the board. Only after we dramatically increase the number of mixtures does the continuous HMM start to offer improved recognition accuracy. The SCHMM can typically reduce the discrete HMM error rate by 10–15% across the board. The continuous HMM with 20 diagonal Gaussian density functions performed worse than either the discrete or the SCHMM when the size of training data was small. It outperformed either the discrete HMM or the SCHMM when sufficient amounts of training data became available. When the amount of training data is sufficiently large, it can reduce the error rate of the semicontinuous HMM by 15–20%.

#### 9.5.2. Isolated vs. Continuous Speech Training

If we build a word HMM for each word in the vocabulary for isolated speech recognition, the training or recognition can be implemented directly, using the basic algorithms introduced in Chapter 8. To estimate model parameters, examples of each word in the vocabulary are collected. The model parameters are estimated from all these examples using the forward-backward algorithm and the reestimation formula. It is not necessary to have precise end-point detection, because the silence model automatically determines the boundary if we concatenate silence models with the word model in both ends.

If subword units,<sup>3</sup> such as phonetic models, are used, we need to share them across different words for large-vocabulary speech recognition. These subword units are concatenated to form a word model, possibly adding silence models at the beginning and end, as illustrated in Figure 9.9.

To concatenate subword units to form a word model, you can have a null transition from the final state of the previous subword HMM to the initial state of the next subword HMM, as indicated by the dotted line in Figure 9.9. As described in Chapter 8, you can estimate the parameters of the concatenated HMM accordingly. Please notice that the added null transition arc should satisfy the probability constraint with the transition probability of each phonetic HMM. The self-loop transition probability of the last state in each individual HMM has the topology illustrated in Figure 9.9. If we estimate these parameters with the concatenated model, the null arc transition probability,  $a_{ij}^{e}$ , should satisfy the constraint  $\sum_{i} (a_{ij} + a_{ij}^{e}) = 1$  such that the self-loop transition probability of the last state is no longer equal to 1. For interword concatenation or concatenation involving multiple pronunciations, you can use multiple null arcs to concatenate individual models together.

We have a detailed discussion on word models vs. subword models in Section 9.4.1.





In the example given in Figure 9.9, we have ten English digits in the vocabulary. We build an HMM for each English phone. The dictionary provides information on each word's pronunciation. We have a special word, *Silence*, that maps to a */sil/* HMM that has the same topology as the standard phonetic HMM. For each word in the vocabulary we first derive the phonetic sequence for each word from the dictionary. We link these phonetic models together to form a word HMM for each word in the vocabulary. The link between two phonetic models is shown in the figure as the dotted arrow.

For example, for word *two*, we create a word model based on the beginning silence *lsill*, phone *ltl*, phone *luwl*, and ending silence *lsill*. The concatenated word model is then treated in the same manner as a standard large composite HMM. We use the standard forward-backward algorithm to estimate the parameters of the composite HMM from multiple sample utterances of the word *two*. After several iterations, we automatically get the HMM parameters for *lsill*, *ltl*, and *luwl*. Since a phone can be shared across different words, the phonetic parameters may be estimated from acoustic data in different words.

The ability to automatically align each individual HMM to the corresponding unsegmented speech observation sequence is one of the most powerful features in the forwardbackward algorithm. When the HMM concatenation method is used for continuous speech, you need to compose multiple words to form a sentence HMM based on the transcription of the utterance. In the same manner, the forward-backward algorithm absorbs a range of pos-

442

#### Acoustic Modeling—Scoring Acoustic Features

sible word boundary information of models automatically. There is no need to have a precise segmentation of the continuous speech.

In general, to estimate the parameters of the HMM, each word is instantiated with its concatenated word model (which may be a concatenation of subword models). The words in the sentence are concatenated with optional silence models between them. If there is a need to modify interword pronunciations due to interword pronunciation change, such as *want you*, you can add a different optional phonetic sequence for *t*-*y* in the concatenated sentence HMM.

In the digit recognition example, if we have a continuous training utterance one three, we compose a sentence HMM, as shown in Figure 9.10, where we have an optional silence HMM between the words one and three, linked with a null transition from the last state of the word model one to the first state of the word model three. There is also a direct null arc connection between the models one and three because a silence may not exist in the training example. These optional connections ensure that all the possible acoustic realizations of the natural continuous speech are considered, so that the forward-backward algorithm can automatically discover the correct path and accurately estimate the corresponding HMM from the given speech observation.

In general, the concatenated sentence HMM can be trained using the forwardbackward algorithm with the corresponding observation sequence. Since the entire sentence HMM is trained on the entire observation sequence for the corresponding sentence, most possible word boundaries are inherently considered. Parameters of each model are based on those state-to-speech alignments. It does not matter where the word boundaries are. Such a training method allows complete freedom to align the sentence model against the observation, and no explicit effort is needed to find word boundaries.

In speech decoding, a word may begin and end anywhere within a given speech signal. As word boundaries cannot be detected accurately, all possible beginning and end points have to be accounted for. This converts a linear search (as for isolated word recognition) to a tree search, and a polynomial recognition algorithm to an exponential one. How to design an efficient decoder is discussed in Chapters 12 and 13.



Figure 9.10 A composite sentence HMM. Each word can be a word HMM or a composite phonetic word HMM, as illustrated in Figure 9.9.

443

### 9.6. ADAPTIVE TECHNIQUES-MINIMIZING MISMATCHES

As Figure 1.2 illustrated, it is important to adapt both acoustic models and language models for new situations. A decent model can accommodate a wide range of variabilities. However, the mismatch between the model and operating conditions always exists. One of the most important factors in making a speech system usable is to minimize the possible mismatch dynamically with a small amount of *calibration data*. Adaptive techniques can be used to modify system parameters to better match variations in microphone, transmission channel, environment noise, speaker, style, and application contexts. As a concrete example, speaker-dependent systems can provide a significant word error-rate reduction in comparison to speaker-independent systems if a large amount of speaker-dependent training data exists [50]. Speaker-adaptive techniques can bridge the gap between these two configurations with a small fraction of the speaker-specific training data needed to build a full speaker-dependent system. These techniques can also be used incrementally as more speech is available from a particular speaker. When speaker-adaptive models are built, you can have not only improved accuracy but also improved speed and potentially reduced model parameter sizes because of accurate representations, which is particularly appealing for practical speech recognition.

There are a number of ways to use adaptive techniques to minimize mismatches. You can have a nonintrusive adaptation process that works in the background all the time. This is typically unsupervised, using only the outcome of the recognizer (with a high confidence score, as discussed in Section 9.7) to guide the model adaptation. This approach can continuously modify the model parameters so that any nonstationary mismatches can be eliminated. As discussed in Chapter 13, systems that are required to transcribe speech in a non-real-time fashion may use multiple recognition passes. You can use unsupervised adaptation on the test data to improve the models after each pass to improve performance for a subsequent recognition pass.

Since the use of recognition results may be imperfect, there is a possibility of divergence if the recognition error rate is high. If the error rate is low, the adaptation results may still not be as good as supervised adaptation in which the correct transcription is provided for the user to read, a process referred to as the *enrollment process*. In this process you can check a wide range of parameters as follows:

- Check the background noise by asking the user not to speak.
- Adjust the microphone gain by asking the user to speak normally.
- Adapt the acoustic parameters by asking the user to read several sentences.
- Change the decoder parameters for the best speed with no loss of accuracy.
- Compose dynamically new enrollment sentences based on the user-specific error patterns.

The challenge for model adaptation is that we can use only a small amount of observable data to modify model parameters. This constraint requires different modeling strategies

## Adaptive Techniques-Minimizing Mismatches

from the ones we discussed in building the baseline system, as the amount of training data is generally sufficient for offline training. In this section we focus on a number of adaptive techniques that can be applied to compensate either speaker or environment variations. Most of these techniques are model-based, since the acoustic model parameters rather than the acoustic feature vectors are adapted. We use speaker-adaptation examples to illustrate how these techniques can be used to improve system performance. We can generalize to environment adaptation by using environment-specific adaptation data and a noise-compensation model, which we discuss in Chapter 10. In a similar manner, we can modify the language model as discussed in Chapter 11.

#### 9.6.1. Maximum a Posteriori (MAP)

Maximum a posteriori (MAP) estimation, as discussed in Chapter 4, can effectively deal with data-sparse problems, as we can take advantage of prior information about existing models. We can adjust the parameters of pretrained models in such a way that limited new training data would modify the model parameters guided by the prior knowledge to compensate for the adverse effect of a mismatch [35]. The prior density prevents large deviations of the parameters unless the new training data provide strong evidence.

More specifically, we assume that an HMM is characterized by a parameter vector  $\Phi$  that is a random vector, and that prior knowledge about the random vector is available and characterized by a prior probability density function  $p(\Phi)$ , whose parameters are to be determined experimentally.

With the observation data X, the MAP estimate is expressed as follows:

$$\hat{\Phi} = \arg\max[p(\Phi \mid \mathbf{X})] = \arg\max[p(\mathbf{X} \mid \Phi)p(\Phi)]$$
(9.14)

If we have no prior information,  $p(\Phi)$  is the uniform distribution, and the MAP estimate becomes identical to the ML estimate. We can use the EM algorithm as the ML to estimate the parameters of HMMs. The corresponding Q-function can be defined as:

$$Q_{MP}(\Phi, \Phi) = \log p(\Phi) + Q(\Phi, \Phi)$$
(9.15)

The EM algorithm for the ML criterion can be applied here directly. The actual expression depends on the assumptions made about the prior density. For the widely used continuous Gaussian mixture HMM, there is no joint conjugate prior density. We can assume different components of the HMM to be mutually independent, so that the optimization can be split into different subproblems involving only a single component of the parameter set. For example, the prior density function for the mixture Gaussian can be as follows:

$$P_{b_i}(\mathbf{c}_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = p_{c_i}(\mathbf{c}_i) \prod_k p_{b_k}(\boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{ik})$$
(9.16)

where  $p_{c_i}(\mathbf{c}_i)$  is a Dirichlet prior density for the mixing coefficient vector of all mixture components in the Markov state *i*, and  $p_{b_k}(\boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{ik})$  denotes the prior density for parameters

Amazon/VB Assets Exhibit 1012 Page 471

(0.10)

of the kth Gaussian component in the state *i*. The Dirichlet prior density  $p_{c_i}(\mathbf{e}_i)$  is characterized by a vector  $v_i$  of positive hyperparameters such that:

$$p_{c_i}(\mathbf{c}_i) \propto \prod_k C_{ik}^{v_k - i} \tag{9.17}$$

For full covariance *D*-dimensional Gaussian densities, the prior density can be a normal-Wishart density parameterized by two values  $\eta > D-1, \tau > 0$ , the vector  $\mu_{nr}$ , and the symmetric positive definite matrix S as follows:

$$p_{b_{k}}(\mu_{ik}, \Sigma_{ik}) \propto \sqrt{\det(\Sigma_{ik})^{D-\eta}} \exp(-\frac{\tau}{2}(\mu_{ik} - \mu_{me})\Sigma_{ik}^{-1}(\mu_{ik} - \mu_{me})' - \frac{1}{2}tr(S\Sigma_{ik}^{-1}))$$
(9.18)

We can apply the same procedure as the MLE Baum-Welch reestimation algorithm. For example, with the Q-function defined in Eq. (9.15), we can apply the Lagrange method to derive the mixture coefficients as follows:

$$\begin{cases} \frac{\partial}{\partial \hat{c}_{ik}} (\log p_{c_i}(\hat{\mathbf{c}}_i) + \sum_k \sum_i \xi_i(i,k) \log \hat{c}_{ik}) + \lambda = 0, \forall k \\ \sum_k \hat{c}_{ik} = 1 \end{cases}$$
(9.19)

Based on Eqs. (9.17) and (9.19), the solution is:

$$\hat{c}_{ik} = \frac{v_{ik} - 1 + \sum_{i} \xi_i(i,k)}{\sum_{i} (v_{ii} - 1 + \sum_{i} \xi_i(i,l))}$$
(9.20)

A comparison between Eq. (9.20) and the ML estimate Eq. (8.58) shows that the MAP estimate is a weighted average between the mode of the prior density and the ML estimate with proportions given by  $v_{ijk} - 1$  and  $\sum \xi_i(i,k)$ , respectively.

We can optimize Eq. (9.15) with respect to mean and covariance parameters in a similar fashion. For example, the solution of these estimates is:

$$\hat{\mu}_{ik} = \frac{\tau_{ik} \mu_{nw_{ik}} + \sum_{i=1}^{T} \zeta_i(i,k) \mathbf{x}_i}{\tau_{ik} + \sum_{i=1}^{T} \zeta_i(i,k)}$$
(9.21)

$$\hat{\Sigma}_{ik} = \frac{S_{ik} + \tau_{ik} (\hat{\mu}_{ik} - \mu_{nw_{ik}}) (\hat{\mu}_{ik} - \mu_{nw_{ik}})^{\prime} + \sum_{i=1}^{T} \zeta_{i} (i,k) (\mathbf{x} - \hat{\mu}_{ik}) (\mathbf{x} - \hat{\mu}_{ik})^{\prime}}{\eta_{ik} - D + \sum_{i=1}^{T} \zeta_{i} (i,k)}$$
(9.22)

Amazon/VB Assets Exhibit 1012 Page 472

## Adaptive Techniques-Minimizing Mismatches

where  $\tau_{ik}$  is the parameter in the normal-gamma density for the corresponding state *i*.

where  $\tau_{k}$  is the reestimation formula for the Gaussian mean is a weighted sum of the prior mean with the ML mean estimate  $\sum_{i=1}^{T} \zeta_{i}(i,k)\mathbf{x}_{i} / \sum_{i=1}^{T} \zeta_{i}(i,k) \cdot \tau_{ik}$  is a balancing factor between prior mean and the ML mean estimate. When  $\tau_{ik}$  is large, the value of the prior knowledge is small and the value of the mean  $\mu_{mv_{ik}}$  is assumed to have high certainty, leading to the dominance of the final estimate. When the amount of adaptation data increases, the MAP estimate approaches the ML estimate, as the adaptation data overwrite any important prior that may influence the final estimate. Similarly, the covariance estimation formula has the same interpretation of the balance between the prior and new data.

One major limitation of the MAP-based approach is that it requires an accurate initial guess for the prior  $p(\Phi)$ , which is often difficult to obtain. We can use the already trained initial models that embody some characteristics of the original training conditions. A typical way to generate an initial Gaussian prior is to cluster the initial training data based on speaker or environment similarity measures. We can derive a set of models based on the partition, which can be seen as a set of observations drawn from a distribution having  $p(\Phi)$ . We can, thus, estimate the prior based on the sample moments to derive the corresponding prior parameters.

Another major limitation is that the MAP-based approach is a local approach to updating the model parameters. Namely, only model parameters that are observed in the adaptation data can be modified from the prior value. When the system has a large number of free parameters, the adaptation can be very slow. Thus in practice we need to find correlations between the model parameters, so that the unobserved or poorly adapted parameters can be altered [3, 22]. Another possibility is to impose structural information so the model parameters can be shared for improved adaptation speed [96].

The MAP training can be iterative, too, which requires an initial estimate of model parameters. A careful initialization for the Gaussian densities is also very important. Unlike the discrete distributions, there is no such a thing as a uniform density for a total lack of information about the value of the parameters. We need to use the same initialization procedure as discussed in Chapter 8.

For speaker-adaptive speech recognition, it has been experimentally found that  $\tau_k$  can be a fixed constant value for all the Gaussian components across all the dimensions. Thus the MAP HMM can be regarded as an interpolated model between the speaker-independent and speaker-dependent HMM. Both are derived from the standard ML forward-backward algorithm. Experimental performance of MAP training is discussed in Section 9.6.3.

## 9.6.2. Maximum Likelihood Linear Regression (MLLR)

When the continuous HMM is used for acoustic modeling, the most important parameter set to adapt is the output Gaussian density parameters, i.e., the mean vector and the covariance matrix. We can use a set of linear regression transformation functions to map both means

# Acoustic Modeling

and covariances in order to maximize the likelihood of the adaptation data [68]. The maximum likelihood linear regression (MLLR) mapping is consistent with the underlying criterion for building the HMM while keeping the number of free parameters under control. Since the transformation parameters can be estimated from a relatively small amount of adaptation data, it is very effective for rapid adaptation. MLLR has been widely used to obtain adapted models for either a new speaker or a new environment condition.

More specifically, in the mixture Gaussian density functions, the *k*th mean vector  $\mu_{t}$  for each state *i* can be transformed using following equation:

$$\tilde{\boldsymbol{\mu}}_{ik} = \mathbf{A}_c \boldsymbol{\mu}_{ik} + \mathbf{b}_c \tag{9.23}$$

where  $\mathbf{A}_c$  is a regression matrix and  $\mathbf{b}_c$  is an additive bias vector associated with some broad class c, which can be either a broad phone class or a set of tied Markov states. The goal of Eq. (9.23) is to map the mean vector into a new space such that the mismatch can be eliminated. Because the amount of adaptation data is small, we need to make sure the number of broad classes c is small so we have only a small number of free parameters to estimate. Equation (9.23) can be simplified into:

$$\tilde{\boldsymbol{\mu}}_{ik} = \mathbf{W}_c \boldsymbol{\mu}_{ik} \tag{9.24}$$

where  $\mu_{ik}$  is extended as  $[1, \mu'_{ik}]'$  and  $W_c$  is the extended transform,  $[b_c, A_c]$ .

This mapping approach is based on the assumption that  $W_c$  can be tied for a wide range of broad phonetic classes so that the overall number of free parameters is significantly less than the number of mean vectors. Therefore, the same transformation can be used for several distributions if they represent similar acoustic characteristics.

To estimate these transformation parameters in the MLE framework, we can use the same Q-function we discussed in Chapter 8. We need to optimize only

$$\sum_{i}\sum_{k=1}^{M} \mathcal{Q}_{b_{i}}(\Phi, \hat{\mathbf{b}}_{ik})$$
(9.25)

with respect to  $\mathbf{W}_c$ . Maximization of  $\mathcal{Q}_b(\Phi, \hat{\mathbf{b}}_{ik})$  with respect to  $\mathbf{W}_c$  can be achieved by computing the partial derivatives. For the Gaussian mixture density function, the partial derivative with respect to  $\mathbf{W}_c$  is:

$$\frac{\partial \hat{b}_{ik}(\mathbf{x})}{\partial \mathbf{W}_{c}} = N(\mathbf{x}, \tilde{\boldsymbol{\mu}}_{ik}, \hat{\boldsymbol{\Sigma}}_{ik}) \boldsymbol{\Sigma}_{ik}^{-1} (\mathbf{x} - \mathbf{W}_{c} \boldsymbol{\mu}_{ik}) \boldsymbol{\mu}_{ik}^{t}$$
(9.26)

Let us denote the set of Gaussian components forming the broad transformation classes as C; we use  $b_{ik} \in C$  to denote that the  $k^{th}$  Gaussian density in state *i* belongs to the

Amazon/VB Assets Exhibit 1012 Page 474

#### Adaptive Techniques-Minimizing Mismatches

class. We can expand the Q-function with the partial derivatives and set it to zero, leading to the following equation:

$$\sum_{i=1}^{T} \sum_{b_{k} \in C} \zeta_{i}(i,k) \Sigma_{ik}^{-1} \mathbf{x}_{i} \boldsymbol{\mu}_{ik}^{\prime} = \sum_{i=1}^{T} \sum_{b_{k} \in C} \zeta_{i}(i,k) \Sigma_{ik}^{-1} \mathbf{W}_{c} \boldsymbol{\mu}_{ik} \boldsymbol{\mu}_{ik}^{\prime}$$
(9.27)

We can rewrite Eq. (9.27) as:

$$\mathbf{Z} = \sum_{\mathbf{b}_{k} \in C} \mathbf{V}_{ik} \mathbf{W}_{c} \mathbf{D}_{ik}$$
(9.28)

where

$$\mathbf{Z} = \sum_{t=1}^{t} \sum_{b_k \in C} \zeta_t(i,k) \Sigma_{ik}^{-1} \mathbf{x}_t \boldsymbol{\mu}_{ik}^t, \qquad (9.29)$$

$$\mathbf{V}_{ik} = \sum_{t=1}^{T} \zeta_t(i,k) \Sigma_{ik}^{-1} , \qquad (9.30)$$

and

$$\mathbf{D}_{ik} = \boldsymbol{\mu}_{ik} \boldsymbol{\mu}_{ik}^{\prime} \,. \tag{9.31}$$

Estimating  $W_c$  for Eq. (9.28) is computationally expensive, as it requires solving simultaneous equations. Nevertheless, if we assume that the covariance matrix is diagonal, we can have a closed-form solution that is computationally efficient. Thus, we can define

$$\mathbf{G}_{q} = \sum_{b_{qk} \in C} v_{qq} \mathbf{D}_{ik} \tag{9.32}$$

where  $v_{qq}$  denotes the  $q^{th}$  diagonal element of matrix  $V_{ik}$ . The transformation matrix can be computed row by row. So for the  $q^{th}$  row of the transformation matrix  $W_q$ , we can derive it from the  $q^{th}$  row of  $Z_q$  [defined in Eq. (9.29)] as follows:

$$\mathbf{W}_q = \mathbf{Z}_q \mathbf{G}_q^{-1} \tag{9.33}$$

Since  $G_q$  may be a singular matrix, we need to make sure we have enough training data for the broad class. Thus, if the amount of training data is limited, we must tie a number of transformation classes together.

We can run several iterations to maximize the likelihood for the given adaptation data. At each iteration, transformation matrices can be initialized to identity transformations. We can iteratively repeat the process to update the means until convergence is achieved. We can

also incrementally adapt the mean vectors after each observation sequence or set of observation sequences while the required statistics are accumulated over time. Under the assumption that the alignment of each observation sequence against the model is reasonably accurate, we can accumulate these estimated counts over time and use them incrementally. In order to deal with the tradeoff between specificity and robust estimation, we can dynamically generate regression classes according to the senone tree. Thus, we can incrementally increase the number of regression classes when more and more data become available.

MLLR adaptation can be generalized to include the variances with the ML framework, although the additional gain after mean transformation is often less significant (less than relative 2% error reduction). When the user donates about 15 sentences for enrollment training, Table 9.5 illustrates how the MLLR adaptation technique can be used to further reduce the word recognition error rate for a typical dictation application. Here, there is only one context-independent phonetic class for all the context-dependent Gaussian densities. As we can see, most of the error reduction came from adapting the mean vectors.

We can further extend MLLR to *speaker-adaptive training* (SAT) [6, 74]. In conventional speaker-independent training, we simply use data from different speakers to build a speaker-independent model. An inherent difficulty in this approach is that spectral variations of different speakers give the speaker-independent acoustic model higher variance than the corresponding speaker-dependent model. We can include MLLR transformation in the process of training to derive the MLLR parameters for each individual speaker. Thus the training data are transformed to maximize the likelihood for the overall speaker-independent model. This process can be run iteratively to reduce mismatches of different speakers. By explicitly accounting for the interspeaker variations during training and decoding, SAT reduces the error rate by an additional 5–10%.

Models	Relative Error Reduction
СНММ	Baseline
MLLR on mean only	+12%
MLLR on mean and variance	+2%
MLLR SA'I	+8%

Table 9.5 Relative error reductions with MLLR methods. The reduction is relative to that of the preceding row.

## 9.6.3. MLLR and MAP Comparison

The MLLR method can be combined with MAP. This guarantees that with the increased amount of training data, we can have not only a set of compact MLLR transformation functions for rapid adaptation, but also directly modified model parameters that converge to the ML estimates. We can use MAP to adapt the model parameters and then add MLLR transform these adapted models. It is also possible to incorporate the MAP principle directly into MLLR [18, 19].

Amazon/VB Assets Exhibit 1012 Page 476

### Adaptive Techniques-Minimizing Mismatches

As an example, the result of a 60,000-word dictation application using various adaptation methods is shown in Figure 9.11.<sup>4</sup> The speaker-dependent model used 1000 utterances. Also included as a reference is the speaker-independent result, which is used as the starting point for adaptive training. When the speaker-independent model is adapted with about 200 utterances, the speaker-adaptive model has already outperformed both speaker-independent and speaker-dependent systems. The results clearly demonstrate that we have insufficient training data for speaker-dependent speech recognition, as MAP-based outperform MLbased models. This also illustrates that we can make effective use of speaker-independent data for speaker-dependent speech recognition. Also, notice that the MLLR method has a faster adaptation rate than the MAP method. The MLLR method has context-independent phonetic classes. So, when the amount of adaptation data is limited, the MLLR method offers better overall performance.

However, the MAP becomes more accurate when the amount of adaptation data increases to 600 per speaker. This is because we can modify all the model parameters with the MAP training, and the MLLR transformation can never have the same degrees of freedom as the MAP method. When the MLLR is combined with MAP, we can have not only rapid adaptation but also superior performance over either the MLLR or MAP method across a wide range of adaptation data points. There are a number of different ways to combine both MLLR and MAP for improved performance [4, 98].



Figure 9.11 Comparison of Whisper with MLLR, MAP, and combined MLLR and MAP. The error rate is shown for a different amount of adaptation data. The speaker-dependent and speaker-independent models are also included. The speaker-dependent model was trained with 1000 sentences.

Amazon/VB Assets Exhibit 1012 Page 477

<sup>&</sup>lt;sup>4</sup> In practice, if a large well-trained, speaker-independent model is used, the baseline performance may be very good, and hence, the relative error reduction from speaker adaptation may be smaller than for smaller and simpler models.

#### 9.6.4. Clustered Models

Both MAP and MLLR techniques are based on using an appropriate initial model for adaptive modeling. How accurate we make the initial model directly affects the overall performance. An effective way to minimize the mismatch is, thus, to cluster similar speakers and environments in the training data, building a set of models for each cluster that has minimal mismatch for different conditions. When we have enough training data, and enough coverage for a wide range of conditions, this approach ensures significantly improved robustness.

For example, we often need a set of clustered models for different telephone channels, including different cellular phone standards. We also need to build gender-dependent models or speaker-clustered models for improved performance. In fact, when we construct speaker-clustered models, we can apply MLLR transformations or neural networks to minimize speaker variations such that different speakers can be mapped to the same golden speaker that is the representative of the cluster.

Speaker clusters can be created based on the information of each speaker-dependent HMM. The clustering procedure is similar to the decision-tree procedure discussed in Section 9.4.3. Using clustered models increases the amount of computational complexity. It also fragments the training data. Clustering is often needed to combine other smoothing techniques, such as deleted interpolation or MLLR transformation, in order to create clustered models from the pooled model. We can also represent a speaker as a weighted sum of individual speaker cluster models with the cluster adaptive training [33] or eigenvoice techniques [64].

When we select an appropriate model, we can compute the likelihood of the test speech against all the models and select the model that has the highest likelihood. Alternatively, we can compute likelihoods as part of the decoding process and prune away less promising models dynamically without significantly increasing the computational load. When multiple models are plausible, we can compute the weighted sum of the clustered models with pretrained mixing coefficients for different clusters, much as we train the deleted interpolation weights.

Traditionally speaker clustering is performed across different speakers without considering phonetic similarities across different speakers. In fact, clustered speaker groups may have very different degrees of variations for different phonetic classes. You can further generalize speaker clustering to the subword or subphonetic level [62]. With multiple instances derived from clustering for each subword or subphonetic unit, you can model speaker variation explicitly across different subword or subphonetic models.

In practice, gender-dependent models can reduce the word recognition error by 10%. More refined speaker-clustered models can further reduce the error rate, but not as much as the gain from gender-dependent models, unless we have a large number of clustered speakers. If the new user happens to be similar to one of these speaker clusters, we can approach speaker-dependent speech recognition without enrollment. For environment-dependent models, clustering is more critical. The challenge is to anticipate the kind of environment or channel distortions the system will have to deal with. Since this is often unpredictable, we

Confidence Measures: Measuring the Reliability

need to use adaptive techniques such as MAP and MLLR to minimize the mismatch. We discuss this in more detail in Chapter 10.

# 9.7. CONFIDENCE MEASURES: MEASURING THE RELIABILITY

One of the most critical components in a practical speech recognition system is a reliable *confidence measure*. With an accurate confidence measure for each recognized word, the conversational back end can repair potential speech recognition errors, can reject out-of-vocabulary words, and can identify key words (perform word spotting) that are relevant to the back end. In a speaker-dependent or speaker-adaptive system, the confidence measure can help user enrollment (to eliminate mispronounced words). It is also critical for unsupervised speaker adaptation, allowing selective use of recognition results so that transcriptions with lower confidence can be discarded for adaptation.

In theory, an accurate estimate of P(W | X), the posterior probability, is itself a good confidence measure for word W given the acoustic input X. Most practical speech recognition systems simply ignore P(X), as it is a constant in evaluating P(W)P(X | W)/P(X) across different words. P(W | X) can be expressed:

$$P(\mathbf{W} \mid \mathbf{X}) = \frac{P(\mathbf{W})P(\mathbf{X} \mid \mathbf{W})}{P(\mathbf{X})} = \frac{P(\mathbf{W})P(\mathbf{X} \mid \mathbf{W})}{\sum_{\mathbf{W}} P(\mathbf{W})P(\mathbf{X} \mid \mathbf{W})}$$
(9.34)

Equation (9.34) essentially provides a solid framework for measuring confidence levels. It is the ratio between the score for the word hypothesis P(W)P(X|W) and the acoustic probability  $\sum_{W} P(W)P(X|W)$ . In the sections that follow we discuss a number of ways to model and use such a ratio in practical systems.

#### 9.7.1. Filler Models

You can compute P(X) in Eq. (9.34) with a general-purpose recognizer. It should be able to recognize anything such that it can *fill the holes* of the grammar in the normal speech recognizer. The filler model has various forms [7, 63]. One of the most widely used is the so-called all-phone network, in which all the possible phonetic and nonspeech HMMs are connected to each other, and with which any word sequence can be recognized.

In addition to evaluating  $P(\mathbf{W})P(\mathbf{X} | \mathbf{W})$  as needed in normal speech recognition, a separate decoding process is used to evaluate  $\sum_{\mathbf{W}} P(\mathbf{W})P(\mathbf{X} | \mathbf{W})$ . Here W is either a phonetic or a word model. You can also apply phonetic *n*-gram probabilities that are derived from a lexicon targeted for possible new words. The best path from the all-phone network is compared with the best path from the normal decoder. The ratio between the two, as expressed in Eq. (9.34), is used to measure the confidence for either word or phone. In the decoding process (see Chapters 12 and 13), you can accumulate the phonetic ratio derived from Eq. (9.34) on a specific word. If the accumulative  $P(\mathbf{W} | \mathbf{X})$  for the word is less than a predetermined threshold, the word is rejected as either a new word or a nonspeech event.

# Acoustic Modeling

Both context-independent and context-dependent phonetic models can be used for the fully connected network. When context-dependent phonetic models are used, you need to make sure that only correct contextual phonetic connections are made. Although context-dependent models offer significant improvement for speech recognition, the filler phonetic network seems to be insensitive to context-dependency in empirical experiments.

There are word-spotting applications that need to spot just a small number of key words. You can use the filler models described here for word spotting. You can also build antiword models trained with all the data that are not associated with the key words of interest. Empirical experiments indicate that large-vocabulary speech recognition is the most suitable choice for word spotting. You can use a general-purpose *n*-gram (see Chapter 11) to generate recognition results and identify needed key words from the word lattice. This is because a large-vocabulary system provides a better estimate of  $\sum_{w} P(W)P(X | W)$  with a more accurate language model probability. In practice, we don't need to use all the hypotheses to compute  $\sum_{w} P(W)P(X | W)$ . Instead, *n*-best lists and scores [40] can be used to provide an effective estimate of  $\sum_{w} P(W)P(X | W)$ .

#### 9.7.2. Transformation Models

To determine the confidence level for each word, subword confidence information is often helpful. Different phones have different impacts on our perception of words. The weight for each subword confidence score can be optimized from the real training data. If a word w has N phones, we can compute the confidence score of the word as follows:

$$CS(w) = \sum_{i=1}^{N} \wp_i(x_i) / N$$
(9.35)

where CS(w) is the confidence score for word w,  $x_i$  is the confidence score for subword unit *i* in word *w*, and  $\mathcal{D}_i$  is the mapping function that may be tied across a number of subword units. The transformation function can be defined as:

 $\wp_i(x) = ax + b \tag{9.36}$ 

We can use discriminative training to optimize the parameters a and b, respectively. A cost function can be defined as a sigmoid function of CS(w). As shown in Figure 9.12, the optimal transformation parameters vary substantially across different phones. The weight for consonants is also typically larger than that of vowels.

The transformation function can be context dependent. Figure 9.13 illustrates the ROC curve of the context-dependent transformation model in comparison with the corresponding phonetic filler model. The filler model essentially has a uniform weight across all the phones in a given word. The estimated transformation model has 15–40% false-acceptance error reduction at various fixed false-rejection rates. The false-acceptance rate of the transformation model is consistently lower than that of the filler model [63].

Amazon/VB Assets Exhibit 1012 Page 480

Confidence Measures: Measuring the Reliability



Figure 9.12 Transformation parameter a for each context-independent phone class. The weight of consonants is typically larger than that of vowels [63].



Figure 9.13 The ROC curve of phonetic filler models with and without optimal feature transformation [63].

### 9.7.3. Combination Models

In practical systems, there are a number of features you can use to improve the performance of confidence measures. For example, the following features are helpful:

- Word stability ratio from different language model weights (*WrdStabRatio*). This is obtained by applying different language weights to see how stably each word shows up in the recognition *n*-best list.
- Logarithm of the average number of active words around the ending frame of the word (*WrdCntEnd*).
- Acoustic score per frame within the word normalized by the corresponding active senone scores (AscoreSen).
- Logarithm of the average number of active words within the word (WrdCntW).
- Acoustic score per frame within the word normalized by the phonetic filler model (AscoreFiller).
- Language model score (LMScore).
- Language model back-off (trigram, bigram, or unigram hit) for the word (LMBackOff).
- Logarithm of the average number of active states within the word (StateCnt).
- Number of phones in the word (Nphones).
- Logarithm of the average number of active words around the beginning frame of the word (*WrdCntBeg*).
- Whether the word has multiple pronunciations (Mpron).
- Word duration (WordDur).

To clarify each feature's relative importance, Table 9.6 shows its linear correlation coefficient against the correct/incorrect tag for each word in the training set. Word stability ratio (*WrdStabRatio*) has the largest correlation value.

Several kinds of classifiers can be used to compute the confidence scores. Previous research has shown that the difference between classifiers, such as linear classifiers, generalized linear models, decision trees, and neural networks, is insignificant. The simplest linear classifier based on discriminative training performs well in practice. As some features are highly correlated, you can iteratively remove features to combat the curse of dimensionality. The combination model can have up to 40–80% false-acceptance error reduction at fixed false-rejection rate in comparison to the single-feature approach.

456

#### Other Techniques

Feature	Correlation
WrdStabRatio	0.590
WrdCntW	-0.223
LMBackOff	0.171
AscoreSen	0.250
LMScore	0.175
Nphones	0.091
WordDur	0.012
WrdCntEnd	0.321
AscoreFiller	0.219
StateCnt	-0.155
Mpron	0.057
WrdCntBeg	-0.067

Table 9.6 Correlation coefficients of several features against correct/incorrect tag.

### 9.8. OTHER TECHNIQUES

In addition to HMMs, a number of interesting alternative techniques are being actively investigated by researchers. We briefly review two promising methods here.

#### 9.8.1. Neural Networks

You have seen both single-layer and multilayer neural nets in Chapter 4 for dealing with static patterns. In dealing with nonstationary signals, you need to address how to map an input sequence properly to an output sequence when two sequences are not synchronous, which should include proper alignment, segmentation, and classification. The basic neural networks are not well equipped to address these problems in a unified way.

Recurrent neural networks have an internal state that is a function of the current input and the previous internal state. A number of them use time-step delayed recurrent loops on the hidden or output units of a feedforward network, as discussed in earlier chapters. For sequences of finite numbers of delays, we can transform these networks into equivalent feedforward networks by unfolding them over the time period. They can be trained with the standard back propagation procedure, with the following modifications:

- The desired outputs are functions of time, and error functions have to be computed for every copy of the output layer. This requires the selection of an appropriate time-dependent target function, which is often difficult to define.
- All copies of the unfolded weights are constrained to be identical during the training. We can compute the correction terms separately for each weight and use the average to update the final estimate.

. .

# Acoustic Modeling

In most of these networks, you can have a partially recurrent network that has feedback of the hidden and output units to the input layer. For example, the feedforward network can be used in a set of local feedforward connections with one time-step delay. These networks are usually implemented by extending the input field with additional feedback units containing both the hidden and output values generated by the preceding input. You can encode the past nonstationary information that is often required to generate the correct output, given the current input, as illustrated in Figure 9.14.



Figure 9.14 A recurrent network with contextual inputs, hidden vector feedback, and output vector feedback.

One of the popular neural networks is the *Time Delay Neural Network* (TDNN) [105]. Like static networks, the TDNN can be trained to recognize a sequence of predefined length (defined by the width of the input window). The activation in the hidden layer is computed from the current and multiple time-delayed values of the preceding layer, and the output units are activated only when a complete speech segment has been processed. A typical TDNN is illustrated in Figure 9.15. The TDNN has been successfully used to classify presegmented phonemes.

All neural networks have been shown to yield good performance for small-vocabulary speech recognition. Sometimes they are better than HMMs for short, isolated speech units. By recurrence and the use of temporal memory, they can perform some kind of integration over time. It remains a challenge for neural networks to demonstrate that they can be as effective as HMMs for dealing with nonstationary signals, as is often required for large-vocabulary speech recognition.

To deal with continuous speech, the most effective solution is to integrate neural nets with HMMs [91, 113]. The neural network can be used as the output probabilities to replace the Gaussian mixture densities. Comparable results can be obtained with the integrated approach. These HMM output probabilities could be estimated by applying the Bayes' rule to the output of neural networks that have been trained to classify HMM state categories. The neural networks can consist either of a single large trained network or of a group of separately trained small networks [21, 31, 75, 90].

> Amazon/VB Assets Exhibit 1012 Page 484

#### Other Techniques

A number of techniques have been developed to improve the performance of training these networks. Training can be embedded in an EM-style process. For example, dynamic programming can be used to segment the training data. The segmented data are then used to retrain the network. It is also possible to have Baum-Welch style training [14, 42].



Figure 9.15 A time-delay neural network (TDNN), where the box  $h_t$  denotes the hidden vector at time t, the box  $x_t$  denotes the input vector at time t, and the box  $z^{-1}$  denotes a delay of one sample.

#### 9.8.2. Segment Models

As discussed in Chapter 8, the HMM *output-independence assumption* results in a piecewise stationary process within an HMM state. Although the nonstationary speech may be modeled sufficiently with a large number of states, the states in which salient speech features are present are far from stationary [25, 99]. While the use of time-derivative features (e.g., delta and/or delta-delta features) alleviates these limitations, the use of such longer-time-span features may invalidate the conditional independence assumption.



Figure 9.16 Diagram illustrating that HMM's output observation can hop between two unexpected quasi-stationary states [46].

The use of Gaussian mixtures for continuous or semicontinuous HMMs, as described in Chapter 8, could introduce another potential problem, where arbitrary transitions among the Gaussian mixture components between adjacent HMM states are allowed [59]. Figure 9.16 illustrates two HMM states with two mixture components. The solid lines denote the valid trajectories actually observed in the training data. However, in modeling these two trajectories, the Gaussian mixtures inadvertently allow two *phantom* trajectories, shown in

# Acoustic Modeling

dashed lines in Figure 9.16, because no constraint is imposed on the mixture transitions across the state. It is possible that such phantom trajectories degrade recognition performance, because the models can overrepresent speech signals that should be modeled by other acoustic units. Segment models can alleviate such HMM modeling deficiencies [77, 79].

In the standard HMM, the output probability distribution is modeled by a quasistationary process, i.e.,

$$P(\mathbf{x}_{1}^{L} \mid s) = \prod_{i=1}^{L} b_{i}(\mathbf{x}_{i})$$
(9.37)

For the segment model (SM), the output observation distribution is modeled by two stochastic processes:

$$P(\mathbf{x}_{1}^{L} \mid s) = P(\mathbf{x}_{1}^{L} \mid s, L)P(L \mid s)$$
(9.38)

The first term of Eq. (9.38) is no longer decomposable in the absence of the outputindependence assumption. The second term is similar to the duration model described in Chapter 8. In contrast to the HMM whose quasi-stationary process for each state *s* generates one frame  $\mathbf{x}_i$ , a state in a segment model can generate a variable-length observation sequence  $\{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_L\}$  with random length *L*.

Since the likelihood evaluation of segment models cannot be decomposed, the computation of the evaluation is not shareable between different segments (even for the case where two segments differ only by one frame). This results in a significant increase in computation for both training and decoding [77]. In general, the search state space is increased by a factor of  $L_{max}$ , the maximum segment duration. If segment models are used for phone segments,  $L_{max}$  could be as large as 60. On top of this large increase in search state space, the evaluation of segment models is usually an order of magnitude more expensive than for HMM, since the evaluation involves several frames. Thus, the segment model is often implemented in a multipass search framework, as described in Chapter 13.

Segment models have produced encouraging performance for small-vocabulary or isolated recognition tasks [25, 44, 79]. However, their effectiveness on large-vocabulary continuous speech recognition remains an open issue because of necessary compromises to reduce the complexity of implementation.

### 9.8.2.1. Parametric Trajectory Models

Parametric trajectory models [25, 37] were first proposed to model a speech segment with curve-fitting parameters. They approximate the *D*-dimensional acoustic observation vector  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$  by a polynomial function. Specifically, the observation vector  $\mathbf{x}_t$  can be represented as

460

Other Techniques

$$\mathbf{x}_{t} = \mathbf{C} \times F_{t} + \mathbf{e}_{t}(\Sigma) = \begin{pmatrix} c_{1}^{0} & c_{1}^{1} & \cdots & c_{1}^{N} \\ c_{2}^{0} & c_{2}^{1} & \cdots & c_{2}^{N} \\ \vdots & \vdots & \vdots & \vdots \\ c_{D}^{0} & c_{D}^{1} & c_{D}^{2} & c_{D}^{N} \end{pmatrix} \begin{pmatrix} f_{0}(t) \\ f_{1}(t) \\ \vdots \\ f_{N}(t) \end{pmatrix} + \mathbf{e}_{t}(\Sigma)$$
(9.39)

where the matrix C is the trajectory parameter matrix,  $F_i$  is the family of  $N^h$ -order polynomial functions, and  $e_i(\Sigma)$  is the residual fitting error. Equation (9.39) can be regarded as modeling the time-varying mean in the output distribution for an HMM state. To simplify computation, the distribution of the residual error is often assumed to be an independent and identically distributed random process with a normal distribution  $N(0, \Sigma)$ . To accommodate diverse durations for the same segment, the relative linear time sampling of the fixed trajectory is assumed [37].

Each segment M is characterized by a trajectory parameter matrix  $C_m$  and covariance matrix  $\Sigma_m$ . The likelihood for each frame can be specified [46] as

$$P(\mathbf{x}_{t} | \mathbf{C}_{m}, \Sigma_{m}) = \frac{\exp\left\{-tr\left[(\mathbf{x}_{t} - \mathbf{C}_{m}F_{t})\Sigma_{m}^{-1}(\mathbf{x}_{t} - \mathbf{C}_{m}F_{t})^{t}\right]/2\right\}}{(2\pi)^{\frac{D}{2}} |\Sigma_{m}|^{\frac{1}{2}}}$$
(9.40)

If we let  $\mathbf{F}_{\mathbf{r}} = (F_0, F_1, \dots, F_{\mathbf{r}-1})^t$ , then the likelihood for the whole acoustic observation vector can be expressed as

$$P(\mathbf{X} | \mathbf{C}_{m}, \Sigma_{m}) = \frac{\exp\left\{-tr\left[(\mathbf{X} - \mathbf{C}_{m}\mathbf{F}_{T})\Sigma_{m}^{-1}(\mathbf{X} - \mathbf{C}_{m}\mathbf{F}_{T})^{t}\right]/2\right\}}{(2\pi)^{\frac{DT}{2}} |\Sigma_{m}|^{\frac{T}{2}}}$$
(9.41)

Multiple mixtures can also be applied to SM. Suppose segment M is modeled by K trajectory mixtures. The likelihood for the acoustic observation vector X becomes

$$\sum_{k=1}^{k} w_k p_k(\mathbf{X} | \mathbf{C}_k, \boldsymbol{\Sigma}_k)$$
(9.42)

Hon et al. [47] showed that only a handful of target trajectories are needed for speakerindependent recognition, in contrast to the many mixtures required for continuous Gaussian HMMs. This should support the phantom-trajectory argument involved in Figure 9.16.

The estimation of segment parameters can be accomplished by the EM algorithm described in Chapter 4. Assume a sample of L observation segments  $X_1, X_2, \dots, X_L$ , with corresponding duration  $T_1, T_2, \dots, T_L$ , are generated by the segment model M. The MLE formulae using the EM algorithm are:

$$\gamma_k^i = \frac{w_k p_k(\mathbf{X}_i | \mathbf{C}_k, \boldsymbol{\Sigma}_k)}{P(\mathbf{X}_i | \mathbf{\Phi}_m)}$$
(9.43)

Amazon/VB Assets Exhibit 1012 Page 487

# Acoustic Modeling

$$\hat{w}_k = \frac{1}{L} \sum_{i=1}^{L} \frac{w_k p_k(\mathbf{X}_i \mid \mathbf{C}_k, \boldsymbol{\Sigma}_k)}{P(\mathbf{X}_i \mid \mathbf{\Phi}_m)}$$
(9.44)

$$\hat{\mathbf{C}}_{k} = \left[\sum_{i=1}^{L} \gamma_{k}^{i} \mathbf{X}_{i} \mathbf{F}_{T_{i}}^{\prime}\right] / \left[\sum_{i=1}^{L} \gamma_{k}^{i} \mathbf{F}_{T_{i}} \mathbf{F}_{T_{i}}^{\prime}\right]$$
(9.45)

$$\hat{\Sigma}_{k} = \sum_{i=1}^{L} \gamma_{k}^{i} (\mathbf{X}_{i} - \mathbf{C}_{k} \mathbf{F}_{T_{i}}) (\mathbf{X}_{i} - \mathbf{C}_{k} \mathbf{F}_{T_{i}})^{\prime} / \sum_{i=1}^{L} \gamma_{k}^{i} T_{i}$$
(9.46)

Parametric trajectory models have been successfully applied to phone classification [25, 46] and word spotting [37], and offer a modestly improved performance over HMMs.

#### 9.8.2.2. Unified Frame- and Segment-Based Models

The strengths of the HMM and the segment-model approaches are complementary. HMMs are very effective in modeling the subtle details of speech signals by using one state for each quasi-stationary region. However, the transitions between quasi-stationary regions are largely neglected by HMMs because of their short durations. In contrast, segment models are powerful in modeling the transitions and longer-range speech dynamics, but might need to give up the detailed modeling to assure trainability and tractability. It is possible to have a unified framework to combine both methods [47].

In the unified complementary framework, the acoustic model  $p(\mathbf{X}_{i}^{T} | \mathbf{W})$  can be considered as two joint hidden processes, as in the following equation:

$$p(\mathbf{X} \mid \mathbf{W}) = \sum_{\mathbf{q}^{h}} \sum_{\mathbf{q}^{i}} p(\mathbf{X}, \mathbf{q}^{h}, \mathbf{q}^{s} \mid \mathbf{W})$$
  
$$= \sum_{\mathbf{q}^{h}} \sum_{\mathbf{q}^{i}} p(\mathbf{X} \mid \mathbf{q}^{h}, \mathbf{q}^{s}) p(\mathbf{q}^{s} \mid \mathbf{q}^{h}) p(\mathbf{q}^{h} \mid \mathbf{W})$$
(9.47)

where  $q^{h}$  represents the hidden process of the HMM and  $q^{s}$ , the segment model. The conditional probability of the acoustic signal  $p(X|q^{s},q^{h})$  can be further decomposed into two separate terms:

$$p(\mathbf{X} \mid \mathbf{q}^{s}, \mathbf{q}^{h}) = p(\mathbf{X} \mid \mathbf{q}^{h}) p(\mathbf{X} \mid \mathbf{q}^{s})^{a}$$
(9.48)

where a is a constant that is called *segment-model weight*. The first term is the contribution from normal frame-based HMM evaluation. We further assume for the second term that recognition of segment units can be performed by detecting and decoding a sequence of salient events in the acoustic stream that are statistically independent. In other words,

$$p(\mathbf{X} \mid \mathbf{q}^s) = \prod_i p(\mathbf{X}_i \mid \mathbf{q}_i^s)$$
(9.49)

where  $\mathbf{X}_i$  denotes the  $i^{th}$  segment.

Amazon/VB Assets Exhibit 1012 Page 488

. ....

#### Other Techniques

We assume that the phone sequence and the phone boundaries hypothesized by HMMs and segment models agree with each other. Based on the independent-segment assumption, this leads to a segment duration model as

$$p(\mathbf{q}^{s} | \mathbf{q}^{h}) = \prod_{i} p(t_{i}, t_{i+1} - 1 | \mathbf{X}_{i})$$
(9.50)

By treating the combination as a hidden-data problem, we can apply the decoding and iterative EM reestimation techniques here. This unified framework enables both frame- and segment-based models to *jointly* contribute to optimal segmentations, which leads to more efficient pruning during the search. The inclusion of the segment models does not require massive revisions in the decoder, because the segment model scores can be handled in the same manner as the language model scores; whereas the segment evaluation is performed at each segment boundary.

Since subphonetic units are often used in HMMs to model the detailed quasistationary speech region, the segment units should be used to model long-range transition. As studies have shown that phone transitions play an essential role in humans' perception, the phone-pair segment unit that spans over two adjacent phones can be used [47]. Let  $e_i$ and  $t_i$  denote the phone and the starting time of the  $i^{th}$  segment, respectively. For a phone-pair  $(e_{i,v}, e_i)$  segment between  $t_i$  and  $t_{i+1}$ , the segment likelihood can be computed as follows:

$$p(\mathbf{X}_{i} \mid q_{i}^{*}) = p(\mathbf{x}_{i-1}^{*} \mid e_{i-1}, e_{i})$$
(9.51)

Rather than applying segment evaluation for every two phones, an *overlapped evalua*tion scheme can be used, as shown in Figure 9.17 (a), where a phone-pair segment model evaluation is applied at each phone boundary. The overlapped evaluation implies that each phone is evaluated twice in the overall score. Most importantly, the overlapped evaluation places constraints on overlapped regions to assure consistent trajectory transitions. This is an important feature, as trajectory mixtures prohibit *phantom* trajectories within a segment unit, but there is still no mechanism to prevent arbitrary trajectory transitions between adjacent segment units.

Some phone-pairs might not have sufficient training data. Units containing silence might also have obscure trajectories due to the arbitrary duration of silence. As a result, a



Figure 9.17 Overlapped evaluation using (a) a phone-pair segment model, or (b) back-off to two phone units when the phone-pair  $(e_{i,i}, e_i)$  segment model does not exist [47].

phone-pair unit  $(e_{i,i}, e_i)$  can be backed off with two phone units as shown in Figure 9.17 (b). The phone units can be context independent or context dependent [46]. Thus, the back-off segment-model evaluation becomes:

$$p(\mathbf{X}_{i} \mid q_{i}^{s}) = \beta * p(\mathbf{x}_{i,-}^{t} \mid e_{i-1}) p(\mathbf{x}_{i,-}^{t} \mid e_{i})$$
(9.52)

where  $\beta$  is the back-off weight, generally smaller than 1.0. The use of back-off weight has the effect of giving more preference to phone-pair segment models than to two-phone-based back-off segment models.

The phone-pair segment model outperformed the phone-pair HMM by more than 20% in a phone-pair classification experiment [46]. The unified framework achieved about 8% word-error-rate reduction on the WSJ dictation task in comparison to the HMM-based Whisper [47].

#### 9.9. CASE STUDY: WHISPER

Microsoft's Whisper engine offers general-purpose speaker-independent continuous speech recognition [49]. Whisper can be used for command and control, dictation, and conversational applications. Whisper offers many features such as continuous speech recognition, speaker-independence with adaptation, and dynamic vocabulary. Whisper has a unified architecture that can be scaled to meet different application and platform requirements.

The Whisper system uses MFCC representations (see Chapter 6) and both first- and second-order delta MFCC coefficients. Two-mean cepstral normalization discussed in Chapter 10 is used to eliminate channel distortion for improved robustness.

The HMM topology is a three-state left-to-right model for each phone. Senone models discussed in Section 9.4.3 are derived from both inter- and intraword context-dependent phones. The generic shared density function architecture can support either semicontinuous or continuous density hidden Markov models.

The SCHMM has a multiple-feature front end. Independent codebooks are built for the MFCC, first-order delta MFCC, second-order delta MFCC, and power and first and second power, respectively. Deleted interpolation is used to interpolate output distributions of context-dependent and context-independent senones. All codebook means and covariance matrices are reestimated together with the output distributions except the power covariance matrices, which are fixed.

The overall senone models can reduce the error rate significantly in comparison to the triphone or clustered triphone model. The shared Markov state also makes it possible to use continuous-density HMMs efficiently for large-vocabulary speech recognition. When a sufficient amount of training data becomes available, the best performance is obtained with the continuous-density mixture HMM. Each senone has 20 mixtures, albeit such an error reduction came at the cost of significantly increased computational complexity.

We can further generalize sharing to the level of each individual Gaussian probability density function. Each Gaussian function is treated as the basic unit to be shared across any Markov state. At this extreme, there is no need to use senones or shared states any more, and

## Historical Perspective and Further Reading

the shared probability density functions become the acoustic kernels that can be used to form any mixture function for any Markov state with appropriate mixture weights. Parameter sharing is, thus, advanced from a phone unit to a Markov state unit (senones) to a density component unit.

Regarding lexicon modeling, most words have one pronunciation in the lexicon. For words that are not in the dictionary, the LTS conversion is based on the decision tree that is trained from the existing lexicon. For the purpose of efficiency, the dictionary is used to store the most frequently used words. The LTS is only used for new words that need to be added on the fly.

For speaker adaptation, the diagonal variances and means are adapted using the MAP method. Whisper also uses MLLR to modify the mean vectors only. The MLLR classes are phone dependent. The transition probabilities are context independent and they are not modified during the adaptation stage.

The language model used in Whisper can be either the trigram or the context-free grammar. The difference is largely related to the decoder algorithm, as discussed in Chapter 13. The trigram lexicon has the 60,000 most-frequent words extracted from a large text corpus. Word selection is based on both the frequency and the word's part-of-speech information. For example, verbs and the inflected forms have a higher weight than proper nouns in the selection process.

Whisper's overall word recognition error rate for speaker-independent continuous speech recognition is about 7% for the standard DARPA business-news dictation test set. For isolated dictation with similar materials, the error rate is less than 3%. If speaker-dependent data are available, it can further reduce the error rate by 15–30%, with less than 30 minutes' speech from each person. The performance can be obtained real-time on today's PC systems.

#### 9.10. HISTORICAL PERSPECTIVE AND FURTHER READING

The first machine to recognize speech was a commercial toy named Radio Rex manufactured in the 1920s. Fueled by increased computing resources, acoustic-phonetic modeling has progressed significantly since then. Relative word error rates have been reduced by 10% every year, as illustrated in Figure 9.18, thanks to the use of HMMs, the availability of large speech and text corpora, the establishment of standards for performance evaluation, and advances in computer technology. Before the HMM was established as the standard, there were many competing techniques, which can be traced back to the 1950s. Gold and Morgan's book provides an excellent historical perspective [38].

The HMM is powerful in that, with the availability of training data, the parameters of the model can be estimated and adapted automatically to give optimal performance. There are many HMM-based state-of-the-art speech recognition systems [1, 12, 27, 34, 49, 55, 72, 73, 93, 108, 109, 112]. Alternatively, we can first identify speech segments, then classify the segments and use the segment scores to recognize words. This approach has produced competitive recognition performance that is similar to HMM-based systems in several small- to medium-vocabulary tasks [115].

> Amazon/VB Assets Exhibit 1012 Page 491

# Acoustic Modeling

Speech recognition systems attempt to model the sources of variability in several ways. At the level of signal representation, in addition to MFCC, researchers have developed representations that emphasize perceptually important speaker-independent features of the signal, and deemphasize speaker-dependent characteristics [43]. Other methods based on linear discriminant analysis to improve class separability [28, 54] and speaker normalization transformation to minimize speaker variations [51, 67, 86, 106, 107, 114] have achieved limited success. Linear discriminant analysis can be traced back to *Fisher's linear discriminant* [30], which projects a dimensional vector onto a single line that is oriented to maximize class separability. Its extension to speech recognition can be found in [65].

At the level of acoustic-phonetic modeling, we need to provide an accurate distance measure of the input feature vectors against the phonetic or word models from the signalprocessing front end. Before the HMM was used, the most successful acoustic-phonetic model was based on the speech template where the feature vectors are stored as the model and dynamic-programming-based time warping was used to measure the distance between the input feature vectors and the word or phonetic templates [88, 94]. The biggest problem for template-based systems is that they are not as trainable as HMMs, since it is difficult to generate a template that is as representative as all the speech samples we have for the particular units of interest.



Figure 9.18 History of DARPA speech recognition word-error-rate benchmark evaluation results from 1988 to 1999. There are four major tasks: the Resource Management command and control task (RM C&C, 1000 words), the Air Travel Information System spontaneous speech understanding task (ATIS, 2000 words), the *Wall Street Journal* dictation task (WSJ, 20,000 words), and the Broadcast News Transcription Task (NAB, 60,000 words) [80-84].

Another approach that attracted many researchers is the knowledge-based one that originated from the Artificial Intelligence research community. This approach requires extensive knowledge engineering, which often led to many inconsistent rules. Due to the com-

Amazon/VB Assets Exhibit 1012 Page 492

## Historical Perspective and Further Reading

plexity of speech recognition, rule-based approaches generally cannot match the performance of data-driven approaches such as HMMs, which can automatically extract salient rules from a large amount of training data [105].

Senones are now widely used in many state-of-the-art systems. Word models or allophone models can also be built by concatenation of basic structures made by states, transitions, and distributions such as *fenones* [8, 9] or *senones* [58].

Segment models, as proposed by Roucos and Ostendorf [79, 92], assume that each variable-length segment is mapped to a fixed number of representative frames. The resulting model is very similar to the HMM with a large number of states. Ostendorf published a comprehensive survey paper [77] on segment models. The parametric trajectory segment model was introduced by Deng et al. [25] and Gish et al. [37] independently. Gish's work is very similar to our description in Section 9.8.2.1, which is based on Hon et al. [46, 47], where the evaluation and estimation are more efficient because no individual polynomial fitting is required for likelihood computation. In addition to the phone-pair units described in this chapter, segment models have also been applied to phonetic units [25], diphones [36], and syllables [78]. The dynamic model [24, 26] is probably the most aggressive attempt to impose a global transition constraint on the speech model. It uses the phonetic target theories on unobserved vocal-tract parameters, which are fed to an MLP to produce the observed acoustic data.

Today, it is not uncommon to have tens of thousands of sentences available for system training and testing. These corpora permit researchers to quantify the acoustic cues important for phonetic contrasts and to determine parameters of the recognizers in a statistically meaningful way. While many of these corpora were originally collected under the sponsorship of the U.S. Defense Advanced Research Projects Agency (ARPA) to spur human language technology development among its contractors [82], they have nevertheless gained international acceptance. Recognition of the need for shared resources led to the creation of the Linguistic Data Consortium (LDC)<sup>5</sup> in the United States in 1992 to promote and support the widespread development and sharing of resources for human language technology. The LDC supports various corpus development activities and distributes corpora obtained from a variety of sources. Currently, LDC distributes about twenty different speech corpora including those cited above, comprising many hundreds of hours of speech. The availability of a large body of data in the public domain, coupled with the specification of evaluation standards, has resulted in uniform documentation of test results, thus contributing to greater reliability in monitoring progress.

To further improve the performance of acoustic-phonetic models, we need a robust system so that performance degrades gracefully (rather than catastrophically) as conditions diverge from those under which it was trained. The best approach is likely to have systems continuously adapted to changing conditions (new speakers, microphone, task, etc.). Such adaptation can occur at many levels in systems, subword models, word pronunciations, language models, and so on. We also need to make the system portable, so that we can rapidly design, develop, and deploy systems for new applications. At present, systems tend to suffer

hup://www.cis.upenn.edu/ldc

significant degradation when moved to a new task. In order to retain peak performance, they must be trained on examples specific to the new task, which is time consuming and expensive. In the new task, system users may not know exactly which words are in the system vocabulary. This leads to a certain percentage of out-of-vocabulary words in natural conditions. Currently, systems lack a very robust method of detecting such out-of-vocabulary words. These words often are inaccurately mapped into the words in the system, causing unacceptable errors.

An introduction to all aspects of acoustic modeling can be found in Spoken Dialogues with Computers [76] and Fundamentals of Speech Recognition [87]. A good treatment of HMM-based speech recognition is given in [52, 60, 105]. Bourlard and Morgan's book [15] is a good introduction to speech recognition based on neural networks. There are a range of applications such as predictive networks that estimate each frame's acoustic vector, given the history [69, 104] and nonlinear transformation of observation vectors [13, 53, 101].

You can find tools to build acoustic models from Carnegie Mellon University's speech open source Web site.<sup>6</sup> This site contains the release of CMU's Sphinx acoustic modeling toolkit and documentation. A version of Microsoft's Whisper system can be found in the Microsoft Speech SDK.<sup>7</sup>

#### REFERENCES

- [1] Abrash, V., et al., "Acoustic Adaptation Using Nonlinear Transformations of HMM Parameters" in Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing 1996, Atlanta, pp. 729-732.
- [2] Acero, A., Acoustical and Environmental Robustness in Automatic Speech Recognition, 1993, Boston, Kluwer Academic Publishers.
- [3] Ahadi-Sarkani, S.M., Bayesian and Predictive Techniques for Speaker Adaptation, Ph. D. Thesis, 1996, Cambridge University, .
- [4] Ahn, S., S. Kang, and H. Ko, "Effective Speaker Adaptations for Speaker Verification," *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey, pp. 1081-1084.
- [5] Alleva, F., et al., "Can Continuous Speech Recognizers Handle Isolated Speech?," Speech Communication, 1998, 26, pp. 183-189.
- [6] Anastasakos, T., et al., "A Compact Model for Speaker Adaptive Training," Int. Conf. on Spoken Language Processing, 1996, Philadelphia, pp. 1137-1140.
- [7] Asadi, A., R. Schwartz, and J. Makhoul, "Automatic Modeling for Adding New Words to a Large-Vocabulary Continuous Speech Recognition System," in Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1991, Toronto, pp. 305-308.
- [8] Bahl, L.R., et al., "Multonic Markov Word Models for Large Vocabulary Continuous Speech Recognition," *IEEE Trans. on Speech and Audio Processing*, 1993, 1(3), pp. 334-344.

http://www.speech.cs.cmu.edu/sphinx/

<sup>&</sup>lt;sup>1</sup> http://www.microsoft.com/speech

#### Historical Perspective and Further Reading

- [9] Bahl, L.R., et al., "A Method for the Construction of Acoustic Markov Models for Words," *IEEE Trans. on Speech and Audio Processing*, 1993, 1(4), pp. 443-452.
- [10] Bahl, L.R., et al., "Automatic Phonetic Baseform Determination," Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1991, Toronto, pp. 173-176.
- [11] Bahl, L.R., et al., "Decision Trees for Phonological Rules in Continuous Speech," Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1991, Toronto, Canada, pp. 185-188.
- [12] Bellegarda, J.R., et al., "Experiments Using Data Augmentation for Speaker Adaptation," Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1995, Detroit, pp. 692-695.
- [13] Bengio, Y., et al., "Global Optimization of a Neural Network-Hidden Markov Model Hybrid," IEEE Trans. on Neural Networks, 1992, 3(2), pp. 252-259.
- [14] Bourlard, H., "A New Training Algorithm for Statistical Sequence Recognition with Applications to Transition-Based Speech Recognition," *IEEE Signal Process*ing Letters, 1996, 3, pp. 203-205.
- [15] Bourlard, H. and N. Morgan, Connectionist Speech Recognition A Hybrid Approach, 1994, Boston, MA, Kluwer Academic Publishers.
- [16] Brown, P.F., *The Acoustic-Modeling Problem in Automatic Speech Recognition*, PhD Thesis in Computer Science Department, 1987, Carnegie Mellon University, Pittsburgh, PA.
- [17] Campbell, J., "Speaker Recognition: A Tutorial," Proc. of the IEEE, 1997, 85(9), pp. 1437-1462.
- [18] Chesta, C., O. Siohan, and C.H. Lee, "Maximum a Posteriori Linear Regression for Hidden Markov Model Adaptation," *Eurospeech*, 1999, Budapest, pp. 211-214.
- [19] Chou, W., "Maximum a Posteriori Linear Regression with Elliptically Symmetric Matrix Priors," *Eurospeech*, 1999, Budapest, pp. 1-4.
- [20] Cohen, M., *Phonological Structures for Speech Recognition*, Ph.D. Thesis 1989, University of California, Berkeley.
- [21] Cook, G. and A. Robinson, "Transcribing Broadcast News with the 1997 Abbot System," Int. Conf. on Acoustics, Speech and Signal Processing, 1998, Seattle, WA, pp. 917-920.
- [22] Cox, S., "Predictive Speaker Adaptation in Speech Recognition," Computer Speech and Language, 1995, 9, pp. 1-17.
- [23] Davis, S. and P. Mermelstein, "Comparison of Parametric Representations for Monosyllable Word Recognition in Continuously Spoken Sentences," *IEEE Trans.* on Acoustics, Speech and Signal Processing, 1980, 28(4), pp. 357-366.
- [24] Deng, L., "A Dynamic, Feature-based Approach to the Interface Between Phonology and Phonetics for Speech Modeling and Recognition," Speech Communication, 1998, 24(4), pp. 299-323.
- [25] Deng, L., et al., "Speech Recognition Using Hidden Markov Models with Polynomial Regression Functions as Nonstationary States," *IEEE Trans. on Speech and Audio Processing*, 1994, 2(4), pp. 507-520.

A	-
Acoustic	M
out	1VIOdelin-

	6
[26]	Digalakis, V., Segment-based Stochastic Models of Spectral Dynamics for Con- tinuous Speech Recognition, Ph.D. Thesis in Electrical Computer System Engi- neering, 1992, Boston University.
[27]	Digalakis, V. and H. Murveit, "Genones: Optimizing the Degree of Mixture Tying in a Large Vocabulary Hidden Markov Model Based Speech Recognizer" in <i>Proc.</i> of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1994, Adelaide, Australia, pp. 537-540.
[28]	Doddington, G.R., "Phonetically Sensitive Discriminants for Improved Speech Recognition," Int. Conf. on Acoustics, Speech and Signal Processing, 1989, Glas- gow, Scotland, pp. 556-559.
[29]	Eichner, M. and M. Wolff, "Data-Driven Generation of Pronunciation Dictionaries in the German Verbmobil Project—Discussion of Experimental Results," <i>IEEE Int.</i> <i>Conf. on Acoustics, Speech and Signal Processing</i> , 2000, Istanbul, Turkey, pp. 1687-1690.
[30]	Fisher, R., "The Use of Multiple Measurements in Taxonomic Problems," Annals of Eugenics, 1936, 7(1), pp. 179-188.
[31]	Fritsch, J. and M. Finke, "ACID/HNN: Clustering Hierarchies of Neural Networks for Context-Dependent Connectionist Acoustic Modeling," Int. Conf. on Acoustics, Speech and Signal Processing, 1998, Seattle, WA, pp. 505-508.
[32]	Fukunaga, K., Introduction to Statistical Pattern Recognition, 2nd ed, 1990, Or- lando, FL, Academic Press.
[33]	Gales, M., "Cluster Adaptive Training for Speech Recognition," Int. Conf. on Spo- ken Language Processing, 1998, Sydney, Australia, pp. 1783-1786.
[34]	Gauvain, J.L., L. Lamel, and M. Adda-Decker, "Developments in Continuous Speech Dictation using the ARPA WSJ Task," <i>Proc. of the IEEE Int. Conf. on</i> <i>Acoustics, Speech and Signal Processing</i> , 1995, Detroit, MI, pp. 65-68.
[35]	Gauvain, J.L. and C.H. Lee, "Bayesian Learning of Gaussian Mixture Densities for Hidden Markov Models," <i>Proc. of the DARPA Speech and Natural Language</i> Workshop, 1991, Palo Alto, CA, pp. 272-277.
[36]	Ghitza, O. and M.M. Sondhi, "Hidden Markov Models with Templates as Non- Stationary States: An Application to Speech Recognition," <i>Computer Speech and</i> Language, 1993, 7(2), pp. 101-120
[37]	Gish, H. and K. Ng, "A Segmental Speech Model with Applications to Word Spot- ting," Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1993, Minneapolis, MN, pp. 447, 450
[38]	Gold, B. and N. Morgan, Speech and Audio Signal Processing: Processing and Perception of Speech and M.
[39]	Greenberg, S., D. Ellis, and J. Hollenback, "Insights into Spoken Language Gleaned from Phonetic Transcription of the Switchboard Corpus," Int. Conf. on Spoken Language Press, 24-27.
[40]	Gunawardana, A., H.W. Hon, and L. Jiang, "Word-Based Acoustic Confidence Measures for Large-Vocabulary Speech Recognition," Int. Conf. on Spoken Lan- guage Processing, 1998, Sydney, Australia, pp. 791-794.
	Amazon/VB Assets

470

Exhibit 1012 Page 496

#### Historical Perspective and Further Reading

- [41] Haeb-Umbach, R., "Investigations on Inter-Speaker Variability in the Feature Space," IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1999, Phoenix, AZ.
- [42] Hennebert, J., et al., "Estimation of Global Posteriors and Forward-Backward Training of Hybrid HMM/ANN Systems," Proc. of the Eurospeech Conf., 1997, Rhodes, Greece, pp. 1951-1954.
- [43] Hermansky, H., "Perceptual Linear Predictive (PLP) Analysis of Speech," Journal of the Acoustical Society of America, 1990, 87(4), pp. 1738-1752.
- [44] Holmes, W. and M. Russell, "Probabilistic-Trajectory Segmental HMMs," Computer Speech and Language, 1999, 13, pp. 3-37.
- [45] Hon, H.W. and K.F. Lee, "CMU Robust Vocabulary-Independent Speech Recognition System," Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1991, Toronto, pp. 889-892.
- [46] Hon, H.W. and K. Wang, "Combining Frame and Segment Based Models for Large Vocabulary Continuous Speech Recognition," *IEEE Workshop on Automatic Speech Recognition and Understanding*, 1999, Keystone, CO.
- [47] Hon, H.-W. and K. Wang, "Unified Frame and Segment Based Models for Automatic Speech Recognition," Int. Conf. on Acoustic, Signal and Speech Processing, 2000, Istanbul, Turkey, IEEE, pp. 1017-1020.
- [48] Huang, X., et al., "From Sphinx II to Whisper: Making Speech Recognition Usable" in Automatic Speech and Speaker Recognition, C.H. Lee, F.K. Soong, and K.K. Paliwal, eds. 1996, Norwell, MA, pp. 481-508, Kluwer Academic Publishers.
- [49] Huang, X., et al., "From Sphinx-II to Whisper Make Speech Recognition Usable" in Automatic Speech and Speaker Recognition, C.H. Lee, F.K. Soong, and K.K. Paliwal, eds. 1996, Norwell, MA, Kluwer Academic Publishers.
- [50] Huang, X. and K.-F. Lee, "On Speaker-Independent, Speaker-Dependent, and Speaker-Adaptive Speech Recognition," *IEEE Trans. on Speech and Audio Proc*essing, 1993, 1(2), pp. 150-157.
- [51] Huang, X.D., "Speaker Normalization for Speech Recognition," Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1992, San Francisco pp. 465-468.
- [52] Huang, X.D., Y. Ariki, and M.A. Jack, *Hidden Markov Models for Speech Recog*nition, 1990, Edinburgh, U.K., Edinburgh University Press.
- [53] Huang, X.D., K. Lee, and A. Waibel, "Connectionist Speaker Normalization and its Applications to Speech Recognition," *IEEE Workshop on Neural Networks for Signal Processing*, 1991, New York, pp. 357-366.
- [54] Hunt, M.J., et al., "An Investigation of PLP and IMELDA Acoustic Representations and of Their Potential for Combination," Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1991, Toronto, pp. 881-884.
- [55] Huo, Q., C. Chan, and C.-H. Lee, "On-Line Adaptation of the SCHMM Parameters Based on the Segmental Quasi-Bayes Learning for Speech Recognition," *IEEE Trans. on Speech and Audio Processing*, 1996, 4(2), pp. 141-144.

472	A coustic Modeling
[56]	Hwang, M.Y., X. Huang, and F. Alleva, "Predicting Unseen Triphones with Senones" Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Proc.
[57]	ing, 1993, Minneapolis, pp. 311-314. Hwang, M.Y. and X.D. Huang, "Acoustic Classification of Phonetic Hidde
[37]	Markov Models" in Proc. of Eurospeech 1991.
[58]	Hwang, M.Y. and X.D. Huang, "Shared-Distribution Hidden Markov Models for Speech Recognition," <i>IEEE Trans. on Speech and Audio Processing</i> , 1993, 1(4) pp. 414-420.
[59]	Iyer, R., et al., "Hidden Markov Models for Trajectory Modeling," Int. Conf. on Spoken Language Processing, 1998, Sydney, Australia.
[60]	Jelinek, F., Statistical Methods for Speech Recognition, 1998, Cambridge, MA MIT Press.
[61]	Jiang, L., H.W. Hon, and X. Huang, "Improvements on a Trainable Letter-to-Sound Converter," <i>Proc. of Eurospeech</i> , 1997, Rhodes, Greece, pp. 605-608.
[62]	Jiang, L. and X. Huang, "Subword-Dependent Speaker Clustering for Improved Speech Recognition," Int Conf. on Spoken Language Processing, 2000, Beijing China.
[63]	Jiang, L. and X.D. Huang, "Vocabulary-Independent Word Confidence Measure Using Subword Features," Int. Conf. on Spoken Language Processing, 1998, Syn- dev Australia
[64]	Kuhn, R., et al., "Eigenvoices for Speaker Adaptation," Int. Conf. on Spoken Lan guage Processing, 1998, Sydney, Australia, pp. 1771-1774
[65]	Kumar, N. and A. Andreou, "Heteroscedastic Discriminant Analysis and Reduced Rank HMMs for Improved Speech Recognition," Speech Communication, 1998 26, pp. 283-297.
[66]	Lee, K.F., Large-Vocabulary Speaker-Independent Continuous Speech Recogni tion: The SPHINX System, Ph.D. Thesis in Computer Science Dept. 1988, Came gie Mellon University. Pittsburgh
[67]	Lee, L. and R. Rose, "Speaker Normalization Using Efficient Frequency Warping Procedures," <i>IEEE Int. Conf. on Acoustics, Speech and Signal Processing</i> , 1996 Atlanta, GA, pp. 353-356
[68]	Leggetter, C.J. and P.C. Woodland, "Maximum Likelihood Linear Regression for Speaker Adaptation of Continuous Density Hidden Markov Models," Computer Speech and Language 1005 0 pp. 171-105
[69]	Levin, E., "Word Recognition Using Hidden Control Neural Architecture," Proc. 0, the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1990, Albuquer out NM np. 422-426
[70]	Lippmann, R.P., E.A. Martin, and D.P. Paul, "Multi-Style Training for Robust Iso lated-Word Speech Recognition," Int. Conf. on Acoustics, Speech and Signal Proc.
[71]	Lucassen, J.M. and R.L. Mercer, "An Information-Theoretic Approach to the Automatic Determination of Phonemic Baseforms," <i>Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing</i> , 1984, San Diego, pp. 42.5.1-42.5.4.

#### Historical Perspective and Further Reading

- [72] Lyu, R.Y., et al., "Golden Mandarin (III) A User-Adaptive Prosodic Segment-Based Mandarin Dictation Machine for Chinese Language with Very Large Vocabulary" in Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1995, Detroit, MI, pp. 57-60.
- [73] Matsui, T., T. Matsuoka, and S. Furui, "Smoothed N-best Based Speaker Adaptation for Speech Recognition" in *Proc. of the IEEE Int. Conf. on Acoustics, Speech* and Signal Processing, 1997, Munich, Germany, pp. 1015-1018.
- [74] McDonough, J., et al., "Speaker-Adapted Training on the Switchboard Corpus" in Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing 1997, Munich, Germany, pp. 1059-1062.
- [75] Morgan, N. and H. Bourlard, Continuous Speech Recognition: An Introduction to Hybrid HMM/Connectionist Approach, in IEEE Signal Processing Magazine, 1995, pp. 25-42.
- [76] Mori, R.D., Spoken Dialogues with Computers, 1998, London, Academic Press.
- [77] Ostendorf, M., V.V. Digalakis, and O.A. Kimball, "From HMM's to Segment Models: a Unified View of Stochastic Modeling for Speech Recognition," *IEEE Trans. on Speech and Audio Processing*, 1996, 4(5), pp. 360-378.
- [78] Ostendorf, M. and K. Ross, "A Multi-Level Model for Recognition of Intonation Labels" in *Computing Prosody*, Y. Sagisaka, W.N. Campell, and N. Higuchi, eds., 1997, New York, pp. 291-308, Springer Verlag.
- [79] Ostendorf, M. and S. Roukos, "A Stochastic Segment Model for Phoneme-Based Continuous Speech Recognition," *IEEE Trans. on Acoustics, Speech and Signal Processing*, 1989, 37(1), pp. 1857-1869.
- [80] Pallett, D., J.G. Fiscus, and J.S. Garofolo, "DARPA Resource Management Benchmark Test Results June 1990," in Proc. of the DARPA Speech and Natural Language Workshop 1990, Hidden Valley, PA, pp. 298-305, Pallett.
- [81] Pallett, D., J.G. Fiscus, and J.S. Garofolo, "DARPA Resource Management Benchmark Test Results," Proc. of the DARPA Speech and Natural Language Workshop, 1991, Morgan Kaufmann Publishers, pp. 49-58.
- [82] Pallett, D.S., et al., "The 1994 Benchmark Tests for the ARPA Spoken Language Program" in Proc. of the ARPA Spoken Language Technology Workshop, 1995, Austin, TX, pp. 5-38.
- [83] Pallett, D.S., et al., "1997 Broadcast News Benchmark Test Results: English and Non-English," Proc. of the Broadcast News Transcription and Understanding Workshop, 1998, Landsdowne, Virginia, Morgan Kaufmann Publishers.
- [84] Pallett, D.S., J.G. Fiscus, and M.A. Przybocki, "1996 Preliminary Broadcast News Benchmark Tests," Proc. of the DARPA Speech Recognition Workshop, 1997, Chantilly, VA, Morgan Kaufmann Publishers.
- [85] Pereira, F., M. Riley, and R. Sproat, "Weighted Rational Transductions and Their Application to Human Language Processing," Proc. of the ARPA Human Language Technology Workshop, 1994, Plainsboro, NJ, pp. 249-254.

474	Acoustic Modeling
[86]	Pye, D. and P.C. Woodland. "Experiments in Speaker Normalization and Adapta- tion for Large Vocabulary Speech Recognition," in <i>Proc. of the IEEE Int. Conf. on</i> Acoustics, Speech and Signal Processing 1997, Munich, Germany, pp. 1047-1050
[87]	Rabiner, L.R. and B.H. Juang, Fundamentals of Speech Recognition, May, 1993, Prentice-Hall.
[88]	Rabiner, L.R. and S.E. Levinson, "Isolated and Connected Word Recognition - Theory and Selected Applications," <i>IEEE Trans. on Communication</i> , 1981, COM- <b>29</b> (5), pp. 621-659.
[89]	Riley, M. and A. Ljolje, eds. Automatic Generation of Detailed Pronunciation Lexicons, in Automatic Speech and Speaker Recognition, ed. C. Lee, F. Soong, and K. Paliwal, 1996, Kluwer Academic Publishers.
[90]	Robinson, A., "An Application of Recurrent Nets to Phone Probability Estimation," IEEE Trans. on Neural Networks, 1994, 5, pp. 298-305.
[91]	Robinson, A.J., et al., "A Neural Network Based, Speaker Independent, Large Vo- cabulary," Proc. of the European Conf. on Speech Communication and Technol- ogy, 1999, Berlin pp. 1941-1944.
[92]	Roucos, S., et al., "Stochastic Segment Modeling Using the Estimate-Maximize Algorithm," Int. Conf. on Acoustic, Speech and Signal Processing, 1988, New York, pp. 127-130.
[93]	Sagisaka, Y. and L.S. Lee, "Speech Recognition of Asian Languages" in Proc. IEEE Automatic Speech Recognition Workshop, 1995, Snowbird, UT, pp. 55-57.
[94]	Sakoe, H. and S. Chiba, "Dynamic Programming Algorithm Optimization for Spo- ken Word Recognition," <i>IEEE Trans. on Acoustics, Speech and Signal Processing</i> , 1978, 26(1), pp. 43-49
[95]	Saraclar, M. and S. Khudanpur, "Pronunciation Ambiguity vs. Pronunciation Vari- ability in Speech Recognition," <i>IEEE Int. Conf. on Acoustics, Speech and Signal</i> <i>Processing</i> , 2000 Istanbul Turkey pp. 1679-1682
[96]	Shinoda, K. and C. Lee, "Structural MAP Speaker Adaptation Using Hierarchical Priors," <i>IEEE Workshop on Automatic Speech Recognition and Understanding</i> , 1997, Santa Barbara, CA, pp. 381–388
[97]	Singh, R., B. Raj, and R. Stern, "Automatic Generation of Phone Sets and Lexical Transcriptions," <i>IEEE Int. Conf. on Acoustics, Speech and Signal Processing</i> , 2000. Istanbul Turkey pp. 1601-1604
[98]	Siohan, O., C. Chesta, and C. Lee, "Joint Maximum a Posteriori Estimation of Transformation and Hidden Markov Model Parameters," <i>IEEE Int. Conf. on Acoustics</i> Space and Size 10.
[99]	Siu, M., et al., "Parametric Trajectory Mixtures for LVCSR," Int. Conf. on Spoken
[100]	Sixtus, A., et al., "Recent Improvements of the RWTH Large Vocabulary Speech Recognition System on Spontaneous Speech," IEEE Int. Conf. on Acoustics, Speech and Signal Days
[101]	Sorenson, H., "A Cepstral Noise Reduction Multi-Layer Network," Int. Conf. on Acoustics, Speech and Signal Processing, 1991, Toronto, pp. 933-936.
	Amazon //P. Accata

### Historical Perspective and Further Reading

- [102] Sproat, R. and M. Riley, "Compilation of Weighted Finite-State Transducers from Decision Trees," ACL-96, 1996, Santa Cruz, pp. 215-222.
- [103] Tajchman, G., E. Fosler, and D. Jurafsky, "Building Multiple Pronunciation Models for Novel Words Using Exploratory Computational Phonology," *Eurospeech*, 1995, pp. 2247-2250.
- [104] Tebelskis, J. and A. Waibel, "Large Vocabulary Recognition Using Linked Predictive Neural Networks," Int. Conf. on Acoustics, Speech and Signal Processing, 1990, Albuquerque, NM, pp. 437-440.
- [105] Waibel, A.H. and K.F. Lee, *Readings in Speech Recognition*, 1990, San Mateo, CA, Morgan Kaufman Publishers.
- [106] Watrous, R., "Speaker Normalization and Adaptation Using Second-Order Connectionist Networks," *IEEE Trans. on Neural Networks*, 1994, 4(1), pp. 21-30.
- [107] Welling, L., S. Kanthak, and H. Ney, "Improved Methods for Vocal Tract Normalization," *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1999, Phoenix, AZ.
- [108] Wilpon, J.G., C.H. Lee, and L.R. Rabiner, "Connected Digit Recognition Based on Improved Acoustic Resolution," *Computer Speech and Language*, 1993, 7(1), pp. 15-26.
- [109] Woodland, P.C., et al., "The 1994 HTK Large Vocabulary Speech Recognition System," Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1995, Detroit, pp. 73-76.
- [110] Wooters, C. and A. Stolcke, "Multiple-Pronunciation Lexical Modeling in a Speaker Independent Speech Understanding System," Proc. of the Int. Conf. on Spoken Language Processing, 1994, Yokohama, Japan, pp. 1363-1366.
- [111] Young, S.J. and P.C. Woodland, "The Use of State Tying in Continuous Speech Recognition," Proc. of Eurospeech, 1993, Berlin pp. 2203-2206.
- [112] Zavaliagkos, G., R. Schwartz, and J. Makhoul, "Batch, Incremental and Instantaneous Adaptation Techniques for Speech Recognition," Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing, 1995, Detroit, pp. 676-679.
- [113] Zavaliagkos, G., et al., "A Hybrid Segmental Neural Net/Hidden Markov Model System for Continuous Speech Recognition," *IEEE Trans. on Speech and Audio Processing*, 1994, 2, pp. 151-160.
- [114] Zhan, P. and M. Westphal, "Speaker Normalization Based on Frequency Warping" in Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing 1997, Munich, Germany, pp. 1039-1042.
- [115] Zue, V., et al., "The MIT SUMMIT System: A Progress Report," Proc. of DARPA Speech and Natural Language Workshop, 1989, pp. 179-189.

# CHAPTER 10

## **Environmental Robustness**

A speech recognition system trained in the lab with clean speech may degrade significantly in the real world if the clean speech used in training doesn't match real-world speech. If its accuracy doesn't degrade very much under mismatched conditions, the system is called *robust*. There are several reasons why realworld speech may differ from clean speech; in this chapter we focus on the influence of the *acoustical environment*, defined as the transformations that affect the speech signal from the time it leaves the mouth until it is in digital format.

Chapter 9 discussed a number of variability factors that are critical to speech recognition. Because the acoustical environment is so important to practical systems, we devote this chapter to ways of increasing the environmental robustness, including microphone, echo cancellation, and a number of methods that enhance the speech signal, its spectrum, and the corresponding acoustic model in a speech recognition system.

477

## 10.1. THE ACOUSTICAL ENVIRONMENT

The acoustical environment is defined as the set of transformations that affect the speech signal from the time it leaves the speaker's mouth until it is in digital form. Two main sources of distortion are described here: additive noise and channel distortion. Additive noise, such as a fan running in the background, door slams, or other speakers' speech, is common in our daily life. Channel distortion can be caused by reverberation, the frequency response of a microphone, the presence of an electrical filter in the A/D circuitry, the response of the local loop of a telephone line, a speech codec, etc. Reverberation, caused by reflections of the acoustical wave in walls and other objects, can also dramatically alter the speech signal.

#### 10.1.1. Additive Noise

Additive noise can be stationary or nonstationary. Stationary noise, such as that made by a computer fan or air conditioning, has a power spectral density that does not change over time. Nonstationary noise, caused by door slams, radio, TV, and other speakers' voices, has statistical properties that change over time. A signal captured with a close-talking microphone has little noise and reverberation, even though there may be lip smacks and breathing noise. A microphone that is not close to the speaker's mouth may pick up a lot of noise and/or reverberation.

As described in Chapter 5, a signal x[n] is defined as white noise if its power spectrum is flat,  $S_{xx}(f) = q$ , a condition equivalent to different samples being uncorrelated,  $R_{xx}[n] = q\delta[n]$ . Thus, a white noise signal has to have zero mean. This definition tells us about the second-order moments of the random process, but not about its distribution. Such noise can be generated synthetically by drawing samples from a distribution p(x); thus we could have uniform white noise if p(x) is uniform, or Gaussian white noise if p(x) is Gaussian. While typically subroutines are available that generate uniform white noise, we are often interested in white Gaussian noise, as it resembles better the noise that tends to occur in practice. See Algorithm 10.1 for a method to generate white Gaussian noise. Variable x is normally continuous, but it can also be discrete.

White noise is useful as a conceptual entity, but it seldom occurs in practice. Most of the noise captured by a microphone is *colored*, since its spectrum is not flat. *Pink* noise is a particular type of colored noise that has a low-pass nature, as it has more energy at the low frequencies and rolls off at higher frequencies. The noise generated by a computer fan, an air conditioner, or an automobile engine can be approximated by pink noise. We can synthesize pink noise by filtering white noise with a filter whose magnitude squared equals the desired power spectrum.

A great deal of additive noise is nonstationary, since its statistical properties change over time. In practice, even the noises from a computer, an air conditioning system, or an automobile are not perfectly stationary. Some nonstationary noises, such as keyboard clicks, are caused by physical objects. The speaker can also cause nonstationary noises such as lip

## The Acoustical Environment

smacks and breath noise. The *cocktail party effect* is the phenomenon under which a human listener can focus onto one conversation out of many in a cocktail party. The noise of the conversations that are not focused upon is called *babble* noise. When the nonstationary noise is correlated with a known signal, the adaptive echo-canceling (AEC) techniques of Section 10.3 can be used.

#### ALGORITHM 10.1: WHITE NOISE GENERATION

To generate white noise in a computer, we can first generate a random variable  $\rho$  with a Rayleigh distribution:

$$p_{\rho}(\rho) = \rho e^{-\rho^2/2} \tag{10.1}$$

from another random variable r with a uniform distribution between (0, 1),  $p_r(r) = 1$ , by simply

equating the probability mass 
$$p_{\rho}(\rho) |d\rho| = p_{r}(r) |dr|$$
 so that  $\left| \frac{dr}{d\rho} \right| = \rho e^{-\rho^{2}/2}$ ; with integration,

it results in  $r = e^{-\rho^2/2}$  and the inverse is given by

$$p = \sqrt{-2 \ln r}$$

(10.2)

If *r* is uniform between (0, 1), and  $\rho$  is computed through Eq. (10.2), it follows a Rayleigh distribution as in Eq. (10.1). We can then generate Rayleigh white noise by drawing independent samples from such a distribution.

If we want to generate white Gaussian noise, the method used above does not work, because the integral of the Gaussian distribution does not exist in closed form. However, if  $\rho$  follows a Rayleigh distribution as in Eq. (10.1), obtained using Eq. (10.2) where *r* is uniform between (0, 1), and  $\theta$  is uniformly distributed between (0, 2 $\pi$ ), then the white Gaussian noise can be generated as the following two variables *x* and *y*.

$$x = \rho \cos(\theta) \tag{10.3}$$

 $y = \rho \sin(\theta)$ 

They are independent Gaussian random variables with zero mean and unity variance, since the Jacobian of the transformation is given by

$$J = \begin{vmatrix} \frac{\partial p_x}{\partial \rho} & \frac{\partial p_x}{\partial \theta} \\ \frac{\partial p_y}{\partial \rho} & \frac{\partial p_y}{\partial \theta} \end{vmatrix} = \begin{vmatrix} \cos\theta & -\rho\sin\theta \\ \sin\theta & \rho\cos\theta \end{vmatrix} = \rho$$
(10.4)

and the joint density p(x, y) is given by

$$p(x,y) = \frac{p(\rho,\theta)}{J} = \frac{p(\rho)p(\theta)}{\rho} = \frac{1}{2\pi}e^{-\rho^2/2}$$

$$= \frac{1}{2\pi}e^{-(x^2+y^2)/2} = N(x,0,1)N(y,0,1)$$
(10.5)

The presence of additive noise can sometimes change the way the speaker speaks. The *Lombard effect* [40] is a phenomenon by which a speaker increases his vocal effort in the presence of background noise. When a large amount of noise is present, the speaker tends to shout, which entails not only a higher amplitude, but also often higher pitch, slightly different formants, and a different coloring of the spectrum. It is very difficult to characterize these transformations analytically, but recently some progress has been made [36].

#### 10.1.2. Reverberation

If both the microphone and the speaker are in an *anechoic*<sup>1</sup> chamber or in free space, a microphone picks up only the direct acoustic path. In practice, in addition to the direct acoustic path, there are reflections of walls and other objects in the room. We are well aware of this effect when we are in a large room, which can prevent us from understanding if the reverberation time is too long. Speech recognition systems are much less robust than humans and they start to degrade with shorter reverberation times, such as those present in a normal office environment.

As described in Chapter 2, the signal level at the microphone is inversely proportional to the distance r from the speaker for the direct path. For the kth reflected sound wave, the sound has to travel a larger distance  $r_k$ , so that its level is proportionally lower. This reflection also takes time  $T_k = r_k/c$  to arrive, where c is the speed of sound in air.<sup>2</sup> Moreover, some energy absorption a takes place each time the sound wave hits a surface. The impulse response of such filter looks like

$$h[n] = \sum_{k=0}^{\infty} \frac{\rho_k}{r_k} \delta[n - T_k] = \frac{1}{c} \sum_{k=0}^{\infty} \frac{\rho_k}{T_k} \delta[n - T_k]$$
(10.6)

where  $\rho_k$  is the combined attenuation of the *k*th reflected sound wave due to absorption. Anechoic rooms have  $\rho_k \approx 0$ . In general  $\rho_k$  is a (generally decreasing) function of frequency, so that instead of impulses  $\delta[n]$  in Eq. (10.6), other (low-pass) impulse responses are used.

Often we have available a large amount of speech data recorded with a close-talking microphone, and we would like to use the speech recognition system with a far field microphone. To do that we can filter the clean-speech training database with a filter h[n], so that the filtered speech resembles speech collected with the far field microphone, and then retrain the system. This requires estimating the impulse response h[n] of a room. Alternatively, we can filter the signal from the far field microphone with an inverse filter to make it resemble the signal from the close-talking microphone.

An anechoic chamber is a room that has walls made of special fiberglass or other sound-absorbing materials so that it absorbs all echoes. It is equivalent to being in free space, where there are neither walls nor reflecting surfaces.

<sup>&</sup>lt;sup>2</sup> In air at standard atmospheric pressure and humidity the speed of sound is c = 331.4 + 0.6T (m/s). It varies with different media and different levels of humidity and pressure.

## The Acoustical Environment

One way to estimate the impulse response is to play a white noise signal x[n] through a loudspeaker or artificial mouth; the signal y[n] captured at the microphone is given by

$$y[n] = x[n] * h[n] + v[n]$$
(10.7)

where v[n] is the additive noise present at the microphone. This noise is due to sources such as air conditioning and computer fans and is an obstacle to measuring h[n]. The impulse response can be estimated by minimizing the error over N samples

$$E = \frac{1}{N} \sum_{n=0}^{N-1} \left( y[n] - \sum_{m=0}^{M-1} h[m] x[n-m] \right)^2$$
(10.8)

which, taking the derivative with respect to h[m] and equating to 0, results in our estimate  $\hat{h}[l]$ :

$$\frac{\partial E}{\partial h[l]}\Big|_{h[l]=\hat{h}[l]} = \frac{1}{N} \sum_{n=0}^{N-1} \left( y[n] - \sum_{m=0}^{M-1} \hat{h}[m]x[n-m] \right) x[n-l] \\
= \frac{1}{N} \sum_{n=0}^{N-1} y[n]x[n-l] - \sum_{m=0}^{M-1} \hat{h}[m] \left( \frac{1}{N} \sum_{n=0}^{N-1} x[n-m]x[n-l] \right) \\
= \frac{1}{N} \sum_{n=0}^{N-1} y[n]x[n-l] - \hat{h}[l] - \sum_{m=0}^{M-1} \hat{h}[m] \left( \frac{1}{N} \sum_{n=0}^{N-1} x[n-m]x[n-l] - \delta[m-l] \right) = 0$$
(10.9)

Since we know our white process is ergodic, it follows that we can replace time averages by ensemble averages as  $N \rightarrow \infty$ :

$$\lim_{N \to \infty} \frac{1}{N} \sum_{n=0}^{N-1} x[n-m]x[n-l] = E\left\{x[n-m]x[n-l]\right\} = \delta[m-l]$$
(10.10)

so that we can obtain a reasonable estimate of the impulse response as

$$\hat{h}[l] = \frac{1}{N} \sum_{n=0}^{N-1} y[n] x[n-l]$$
(10.11)

Inserting Eq. (10.7) into Eq. (10.11), we obtain

$$\hat{h}[l] = h[l] + e[l] \tag{10.12}$$

where the estimation error e[n] is given by

$$e[l] = \frac{1}{N} \sum_{n=0}^{N-1} \nu[n] x[n-l] + \sum_{m=0}^{M-1} h[m] \left( \frac{1}{N} \sum_{n=0}^{N-1} x[n-m] x[n-l] - \delta[m-l] \right)$$
(10.13)

If v[n] and x[n] are independent processes, then  $E\{e[l]\}=0$ , since x[n] is zero-mean, so that the estimate of Eq. (10.11) is unbiased. The covariance matrix decreases to 0 as