

# CSE 252B: Computer Vision II

Lecturer: Serge Belongie

Scribe: Jayson Smith

## LECTURE 4

### Planar Scenes and Homography

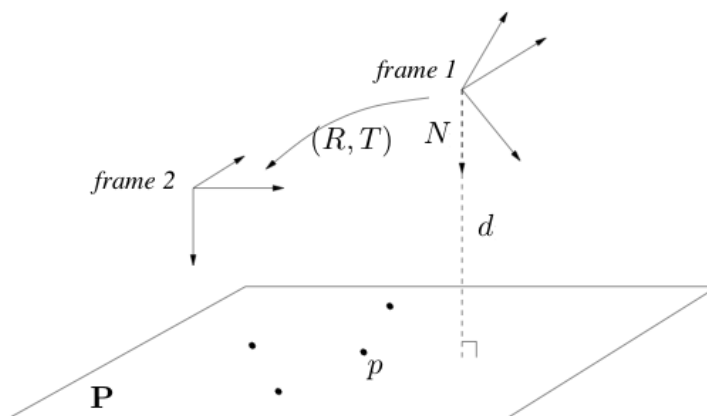
#### 4.1. Points on Planes

This lecture examines the special case of planar scenes. When talking about the 8-point algorithm in two views, we saw that there are configurations of points lying on **critical surfaces** that cause the algorithm to fail. One example is when all points lie on a plane. We need an alternative algorithm to handle this important special case.

---

<sup>1</sup>Department of Computer Science and Engineering, University of California, San Diego.

April 7, 2004



The above figure (from an early draft of MaSKS) is a special case of the two-view configuration of points on a plane:

$$\mathbf{X}_2 = R\mathbf{X}_1 + \mathbf{T}$$

In this case, the vector normal from camera 1 is a perpendicular distance  $d$  from the plane  $P$ :

$$\mathbf{N}^\top \mathbf{X}_1 = n_1X + n_2Y + n_3Z = d$$

or,

$$\frac{1}{d}\mathbf{N}^\top \mathbf{X}_1 = 1 \quad \forall \mathbf{X}_1 \in P$$

## 4.2. Homography

Using the definition of the normal vector above, and multiplying and dividing by  $d$  (which equals  $\mathbf{N}^\top \mathbf{X}_1$ ), the expression for the transformation becomes

$$\mathbf{X}_2 = R\mathbf{X}_1 + \mathbf{T}\frac{1}{d}\mathbf{N}^\top \mathbf{X}_1 = H\mathbf{X}_1$$

where

$$H = R + \frac{1}{d}\mathbf{T}\mathbf{N}^\top, \quad H \in \mathbb{R}^{3 \times 3}$$

$H$  is known as the **planar homography matrix**. Previously, we had only the epipolar constraint which mapped a point in one view to a line in the other view:

$$\mathbf{x}_2^\top E \mathbf{x}_1 = \mathbf{x}_2^\top \hat{T} R \mathbf{x}_1 = 0$$

A homography is stronger: it is a point-to-point mapping from one view to another.

### 4.3. From 3D to 2D Coordinates

Under homography, we can write the transformation of points in 3D from camera 1 to camera 2 as:

$$\mathbf{X}_2 = H\mathbf{X}_1 \quad \mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^3$$

In the image planes, using homogeneous coordinates, we have

$$\lambda_1 \mathbf{x}_1 = \mathbf{X}_1, \quad \lambda_2 \mathbf{x}_2 = \mathbf{X}_2, \quad \text{therefore} \quad \lambda_2 \mathbf{x}_2 = H\lambda_1 \mathbf{x}_1$$

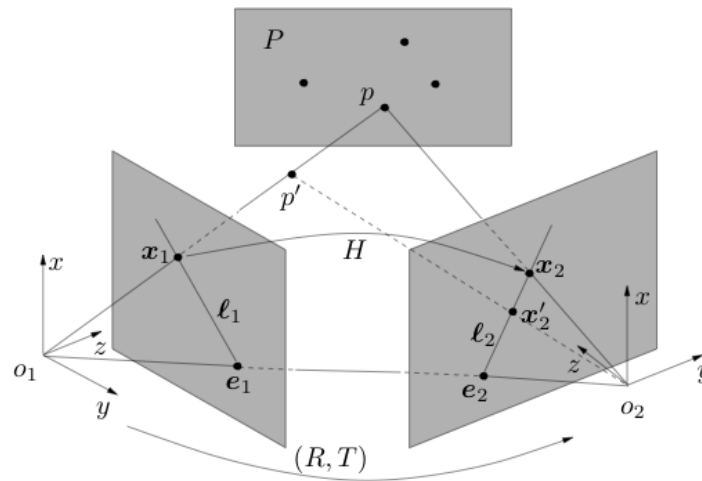
This means that  $\mathbf{x}_2$  is equal to  $H\mathbf{x}_1$  up to a scale (due to universal scale ambiguity). The consequence of this ambiguity is that  $\mathbf{T}$  and  $d$  can't be individually extracted from  $H$ . We can only extract  $\frac{\mathbf{T}}{d}$ .

Note that  $\mathbf{x}_2 \sim H\mathbf{x}_1$  is a *direct mapping* between points in the image planes.

If it is known that some points all lie in a plane in the scene, the image can be rectified directly without needing to recover and manipulate 3D coordinates.

### 4.4. Induced Homography

We say the homography  $H$  is *induced* by coplanar points in the scene.



The pose  $(R, \mathbf{T})$  is a 3D relationship, but  $H$  is a shortcut that directly relates points in the image planes. When the homography is applied to a point  $p'$  that does not actually lie on plane  $P$ , it will get mapped in image plane 2 as if it were “splatted” on  $P$  along the ray  $\overline{o_1 p'}$ .

For a point  $p'$  not on plane  $P$ , the epipolar lines are given by:

$$\mathbf{l}_2 \sim \hat{x}'_2 H \mathbf{x}_1 \quad \text{and} \quad \mathbf{l}_1 \sim H^\top \mathbf{l}_2$$

$\mathbf{x}'_2$  is the image of point  $p'$  in image plane 2, and  $H\mathbf{x}_1$  is where the homography maps  $\mathbf{x}_1$  in image plane 2. These two points lie on  $\mathbf{l}_2$ , and their cross product gives us the  $\mathbf{l}_2$ .

The transformation of lines under the homography  $H$  arises as follows:

$$\begin{aligned} \mathbf{l}_2^\top \mathbf{x}_2 &= 0 \quad \forall \mathbf{x}_2 \\ \Rightarrow \mathbf{l}_2^\top H\mathbf{x}_1 &= 0 \quad \forall \mathbf{x}_1 \quad (\mathbf{l}_2^\top H \text{ must be the epipolar line through } \mathbf{x}_1) \\ \Rightarrow \mathbf{l}_1 &= H^\top \mathbf{l}_2 \\ \Rightarrow \mathbf{l}_2 &= H^{-\top} \mathbf{l}_1 \end{aligned}$$

Thus we have means to relate points and lines from camera 1 to camera 2:

$$\begin{aligned} \mathbf{x}_2 &= H\mathbf{x}_1 \quad (\text{contravariant}) \\ \mathbf{l}_2 &= H^{-\top} \mathbf{l}_1 \quad (\text{covariant}) \end{aligned}$$

The important point here is that we can recover the epipolar lines without knowing the essential matrix  $E$ .

Note: There is a ‘‘Six point algorithm’’ (which is not its official name) that uses 4 points on a plane to induce a homography, and two points not on the plane to find the epipolar lines. This appears in the next lecture.

## 4.5. Homography Estimation

To estimate  $H$ , we start from the equation  $\mathbf{x}_2 \sim H\mathbf{x}_1$  and cross both sides with  $\mathbf{x}_2$ :

$$\begin{aligned} \hat{x}_2 \mathbf{x}_2 &\sim \hat{x}_2 H\mathbf{x}_1 \\ \Rightarrow \hat{x}_2 H\mathbf{x}_1 &= \mathbf{0} \end{aligned}$$

since any vector crossed with itself goes to  $\mathbf{0}$ . This is the *planar homography constraint*.

What would happen if we ignored the fact that these points were on a plane and tried to estimate  $E$ ?

Since  $\mathbf{x}_2 \sim H\mathbf{x}_1$ , then for every  $\mathbf{u} \in \mathbb{R}^3$ , we know:

$$\begin{aligned} \mathbf{u} \times \mathbf{x}_2 &= \hat{\mathbf{u}}\mathbf{x}_2 \perp H\mathbf{x}_1 \\ \Rightarrow (\hat{\mathbf{u}}\mathbf{x}_2)^\top H\mathbf{x}_1 &= 0 \quad (\text{note: } \hat{\mathbf{u}} \text{ is skew-symmetric}) \\ \Rightarrow -\mathbf{x}_2^\top \hat{\mathbf{u}}H\mathbf{x}_1 &= 0 \quad (\text{since homogeneous, can drop minus sign}) \\ \Rightarrow \mathbf{x}_2^\top \hat{\mathbf{u}}H\mathbf{x}_1 &= 0 \\ \Rightarrow \mathbf{x}_2^\top E\mathbf{x}_1 &= 0 \end{aligned}$$

This final equality is true for a *family* of matrices  $E = \hat{\mathbf{u}}H \in \mathbb{R}^{3 \times 3}$  besides the true  $E = \hat{T}R$ . This will cause the 8-point algorithm to crash.

## 4.6. Purely Rotating Camera

A camera that is rotating but has no translation maps points in 3D as

$$\mathbf{X}_2 = R\mathbf{X}_1$$

This is a special case of homography:

$$H = R + \frac{1}{d}\mathbf{T}\mathbf{N}^\top \quad \text{with} \quad \mathbf{T} = \mathbf{0},$$

so

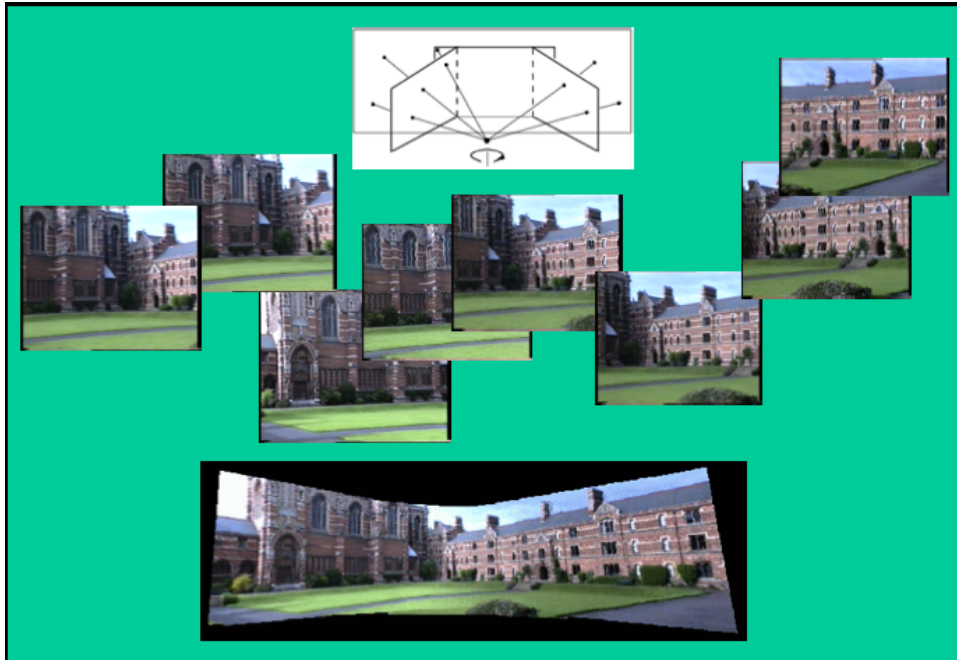
$$\hat{x}_2 H \mathbf{x}_1 = \mathbf{0}$$

becomes

$$\hat{x}_2 R \mathbf{x}_1 = \mathbf{0}$$

A camera rotating about its optical center captures images of a 3D scene as if the scene were painted on a plane infinitely far away from the camera. No depth can be perceived without a translation between the two views. Depth cues (parallax) can only be recovered when  $\mathbf{T}$  is nonzero. Looking at the homography equation, the limit of  $H$  as  $d$  approaches infinity is  $R$ . Thus any pair of images of an arbitrary scene captured by a purely rotating camera is related by a planar homography.

A *planar panorama* can be constructed by capturing many overlapping images at different rotations, picking an image to be a reference, and then finding corresponding points between the overlapping images. The pairwise homographies are derived from the corresponding points, forming a mosaic that typically is shaped like a “bow-tie,” as images farther away from the reference are warped outward to fit the homography. The figure below is from Pollefeys and Hartley & Zisserman.



#### 4.7. Second Derivation of Homography Constraint

The homography constraint, element by element, in homogenous coordinates is as follows:

$$\begin{bmatrix} x_2 \\ y_2 \\ z_2 \end{bmatrix} = \begin{bmatrix} H_{11} & H_{12} & H_{13} \\ H_{21} & H_{22} & H_{23} \\ H_{31} & H_{32} & H_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \\ z_1 \end{bmatrix} \Leftrightarrow \mathbf{x}_2 = H\mathbf{x}_1$$

In inhomogenous coordinates ( $x'_2 = x_2/z_2$  and  $y'_2 = y_2/z_2$ ),

$$x'_2 = \frac{H_{11}x_1 + H_{12}y_1 + H_{13}z_1}{H_{31}x_1 + H_{32}y_1 + H_{33}z_1}$$

$$y'_2 = \frac{H_{21}x_1 + H_{22}y_1 + H_{23}z_1}{H_{31}x_1 + H_{32}y_1 + H_{33}z_1}$$

Without loss of generality, set  $z_1 = 1$  and rearrange:

$$x'_2(H_{31}x_1 + H_{32}y_1 + H_{33}) = H_{11}x_1 + H_{12}y_1 + H_{13}$$

$$y'_2(H_{31}x_1 + H_{32}y_1 + H_{33}) = H_{21}x_1 + H_{22}y_1 + H_{23}$$

We want to solve for  $H$ . Even though these inhomogeneous equations involve the coordinates nonlinearly, the coefficients of  $H$  appear linearly.

Each corresponding point  $\mathbf{x}_1^i \longleftrightarrow \mathbf{x}_2^i$  gives two equations.  $H$  has 8 degrees of freedom, so we need 8 equations = 4 correspondences to get  $H$  (up

to a scale factor). This should match intuition in that one needs 4 corners to describe the mapping of a square under perspective projection.

To solve, stack  $H$  into  $\mathbf{H}^s \in \mathbb{R}^9$  (i.e.  $\mathbf{H}(\cdot)$  in matlab). Write

$$\begin{aligned} \hat{x}_2 H \mathbf{x}_1 & \text{ as} \\ a^\top \mathbf{H}^s &= \mathbf{0}, \quad \text{where} \\ a &\doteq \mathbf{x}_1 \otimes \hat{x}_2 \in \mathbb{R}^{9 \times 3} \end{aligned}$$

Think of this as stuffing all the necessary cross terms between the two sets of coordinates into blocks of a matrix. This allows one to make explicit the linear dependence of the values in  $H$ .

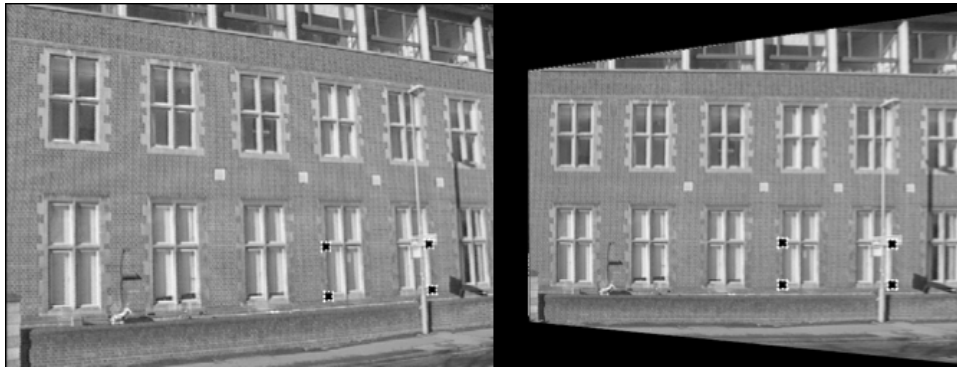
Note: since the rank of a ‘hatted’ matrix is 2, the third equation added by the Kronecker product is a linear combination of one of the other two, and is thus redundant. Also, as convenient as the Kronecker product is here, it does not carry over to the derivation for 3D case, which we’ll need later on in the class.

Collect the  $a$ ’s for each correspondence into a “design matrix”  $\chi$ ,

$$\begin{aligned} \chi &\doteq [a^1 \ a^2 \ \dots \ a^n]^\top \in \mathbb{R}^{3n \times 9}, \quad \text{then} \\ \chi \mathbf{H}^s &= \mathbf{0} \end{aligned}$$

$\text{rank}(\chi) = 8$ ; solve for  $\mathbf{H}^s$  as the null vector, then reshape it into  $H$ .

The process described above is known as the 4-point algorithm or the Direct Linear Transform (DLT). It appears as Algorithm 5.2 in MaSKS.<sup>1</sup> The pair of images below (from Hartley & Zisserman) shows an image before and after projective distortion correction using  $H$  obtained from the four indicated points in the left image and their known coordinates in the scene.



$H$  is homogeneous, specified up to a universal scale factor. In MaSKS Section 5.3.3, they give the decomposition of  $H$  into  $\{R, \frac{1}{d}\mathbf{T}, \mathbf{N}\}$ . There are 4 possible solutions, two of which are physically possible in the sense of positive depth w.r.t. the two image planes.

<sup>1</sup>Here is a question to think about. Why can’t you just solve for  $H$  via a pseudoinverse on  $X_2 = HX_1$ , where  $X_i = [\mathbf{x}_i^1, \dots, \mathbf{x}_i^n]$ ?