



FACULTY OF MATHEMATICS,
PHYSICS AND INFORMATICS

Comenius University
Bratislava

3D Vision

Lecture 3: Vanishing Points and Lines, Homography Estimation

Ing. Viktor Kocur, PhD.

7.3.2023



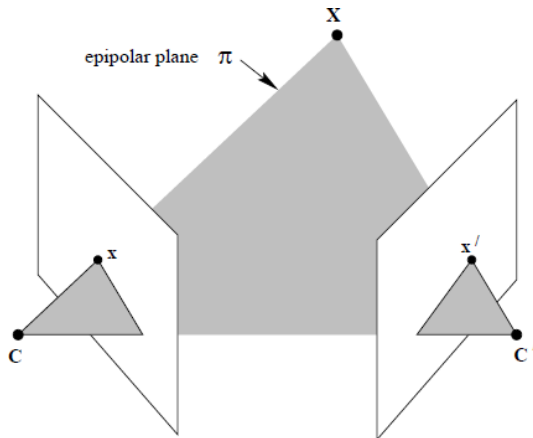
- Fundamental Matrix
- Essential Matrix
- 8-point algorithm
- 7-point algorithm
- 5-point algorithm

How are images in two views related?



Suppose we have two cameras with $P = K[I|\mathbf{0}]$ and $P' = K'[R, \mathbf{t}]$. (Note that we will keep using this notation on future slides.) We want to know how a single point \mathbf{X} from the real scene is imaged in both images as \mathbf{x} and \mathbf{x}' respectively.

What do we know about the relationship of \mathbf{x}, \mathbf{x}' and \mathbf{X} ?



We can see that $\mathbf{x}, \mathbf{x}', \mathbf{X}$ and both of the camera centers are coplanar. We denote this plane as π .



Now suppose that we only know \mathbf{x} . What can we tell about the position of \mathbf{x}' ?

Since the point \mathbf{x}' has to lie on π and also the center of the second camera lies on π too, the point \mathbf{x}' has to lie on line l' which is the image of π .

Epipolar Geometry

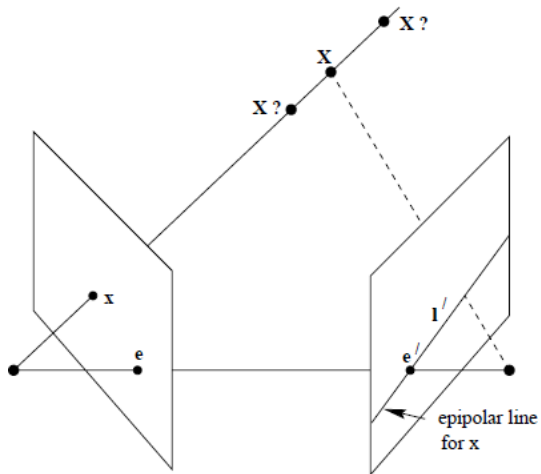


Image adopted from: Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003

Epipolar Plane

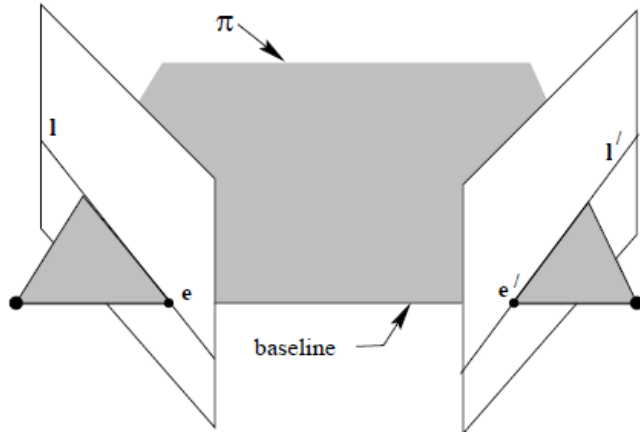


Image adopted from: Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003

Pencil of Epipolar Lines

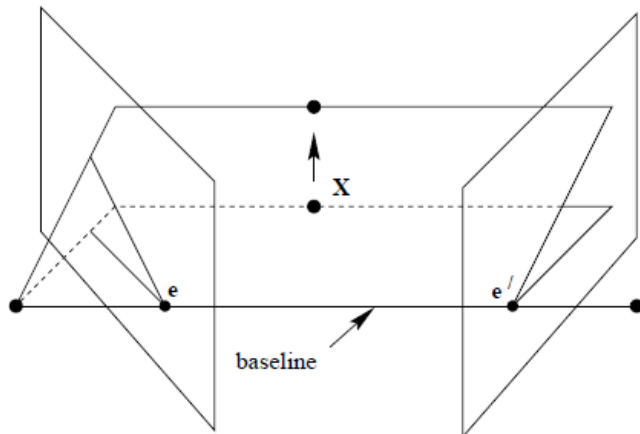


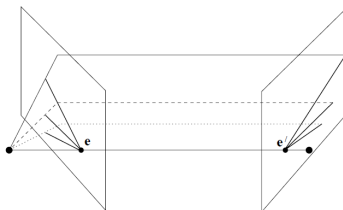
Image adopted from: Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003



Let us now define some of the geometrical objects we used:

- The **epipole** is the point of intersection of line joining the camera centers and the image plane. Equivalently, the epipole is the image in one view of the camera centre of the other view. It is also the vanishing point of the baseline (translation) direction.
- An **epipolar plane** is a plane containing the baseline. There is a one-parameter family (a pencil) of epipolar planes.
- An **epipolar line** is the intersection of an epipolar plane with the image plane. All epipolar lines intersect at the epipole. An epipolar plane intersects the left and right image planes in epipolar lines, and defines the correspondence between the lines.

Pencil of Epipolar Lines



a



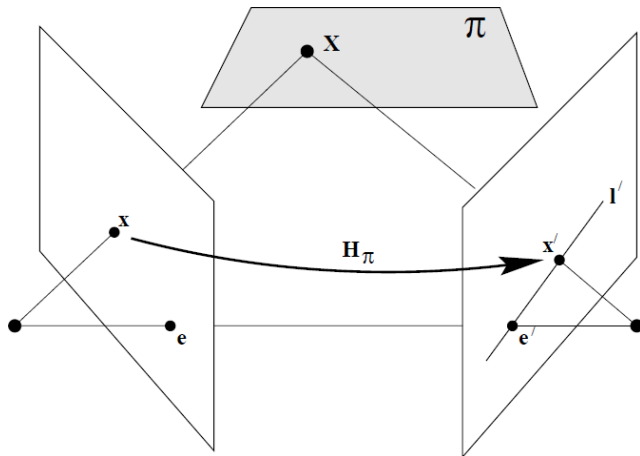
Image adopted from: Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003



We thus know that there is a map $\mathbf{x} \mapsto \mathbf{l}'$. It turns out that this mapping will have the form of:

$$\mathbf{l}' \sim F\mathbf{x}'. \quad (1)$$

The matrix F is the **fundamental matrix**. It encompasses all information about the relative geometry of the two cameras. The rank of F is 2 and it is thus singular.



Consider a plane π which contains \mathbf{X} and does not pass through any of the two camera centers.



In order to obtain the line l' we want to form the line from e' to x' . We can use transfer via plane by defining a homography H_π between the two views. Thus we can get $x' = H_\pi x$. Now, recall that we can construct lines using the vector product of two points and we can also rewrite vector product as matrix multiplication using skew-symmetric matrices. We can thus obtain the final expression:

$$l' = [e']_\times H_\pi x = Fx. \quad (2)$$

This also shows that F is of rank 2, because $[e']_\times$ is rank 2 and H_π is regular. Note that we have freedom in choosing the plane π and also thus H_π . The plane does not have to be present in the scene in reality.



Suppose we have two cameras with $P = K[I|\mathbf{0}]$ and $P' = K'[R, \mathbf{t}]$. Then it is possible to derive the following expressions:

$$F \sim [\mathbf{e}]_{\times} K' R K^{-1} \sim K'^{-T} [\mathbf{t}]_{\times} R K^{-1} \sim K'^{-T} R [\mathbf{R}^T \mathbf{t}]_{\times} K^{-1} \sim K'^{-T} R K^T [\mathbf{e}]_{\times}. \quad (3)$$

Note that F contains information about rotation and translation and both intrinsic camera matrices.



Consider again two cameras P and P' . We can now consider a relationship between two images \mathbf{x} and \mathbf{x}' . Since \mathbf{x}' has to lie on $\mathbf{l}' \sim F\mathbf{x}$ we can write the correspondence condition:

$$\mathbf{x}'^T F \mathbf{x} = 0. \quad (4)$$

Note that this enables us to characterize the F only using points without the need to knowing K, K', R or \mathbf{t} . This is similar to how we defined homographies and then showed that they can be described by 4 correspondences.



- Transpose - if F is the fundamental matrix for pair (P, P') then F^T is the fundamental matrix for (P', P) .
- Epipolar lines - for any point \mathbf{x} in the first image $F\mathbf{x}$ represent the epipolar line \mathbf{l}' in the second image. Equivalently $F^T\mathbf{x}'$ represents \mathbf{l} .
- The epipoles - It holds that $F\mathbf{e} = \mathbf{0}$ and $\mathbf{e}'^T F = \mathbf{0}$.
- F has seven degrees of freedom. It has 9 entries. One degree of freedom is reduced by scaling and another one is reduced due to $\text{rank}(F)$.



We may consider a special case of motion of single camera ($K' = K$) which is pure translation (e.g. $R = I$). This type of motion leads to simplification of F to:

$$F = [\mathbf{e}']_{\times} K K^{-1} = [\mathbf{e}']_{\times}. \quad (5)$$

This leaves F with only two degrees of freedom as we only need to find \mathbf{e}' .

This also gives us insight into F for general motion. We can use a homography H to rotate the second image so that $R = I$ and then we have $F = [\mathbf{e}']_{\times} H$.

Pure Translation

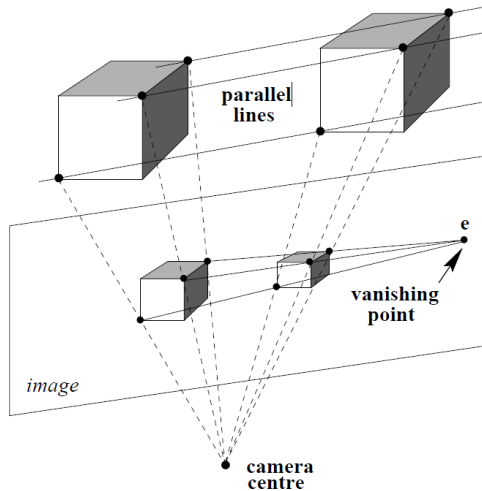


Image adopted from: Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003



Let us now consider a case when the two cameras are calibrated. We can then define $\hat{\mathbf{x}} = K^{-1}\mathbf{x}$ and $\hat{\mathbf{x}}' = K'^{-1}\mathbf{x}'$. This is equivalent to cameras $\hat{P} = [I|\mathbf{0}]$ and $\hat{P}' = [R|\mathbf{t}]$, where $\hat{\mathbf{x}} \sim \hat{P}\mathbf{X}$ and $\hat{\mathbf{x}}' \sim \hat{P}'\mathbf{X}$. From (3) we can see that for a pair of correspondences we have:

$$0 = \mathbf{x}'^T F \mathbf{x} = \hat{\mathbf{x}}'^T K'^T F K \hat{\mathbf{x}} = \hat{\mathbf{x}}'^T K'^T K'^{-T} [\mathbf{t}]_{\times} R K^{-1} K \hat{\mathbf{x}} = \hat{\mathbf{x}}'^T [\mathbf{t}]_{\times} R \hat{\mathbf{x}} = \hat{\mathbf{x}}'^T E \hat{\mathbf{x}} = 0, \quad (6)$$

where matrix $E = [\mathbf{t}]_{\times} R$ is called the **essential matrix**. It relates two views with calibrated cameras. We can see from the derivation above that:

$$E = K'^T F K. \quad (7)$$



E has to be of rank two. However a further condition applies. The singular values of E have to be $(\sigma, \sigma, 0)$. In other words, the two non-zero singular values have to be equal, thus SVD of E is:

$$E = U \text{diag}(\sigma, \sigma, 0) V^T. \quad (8)$$

This leaves E with 5 degrees of freedom as opposed to the 7 degrees for F .

Extracting R and \mathbf{t} from E



We can factorize E to the form $E = [\mathbf{t}]_{\times} R$ where R is a rotation matrix and \mathbf{t} is the translation vector. First let us define two matrices Z and W :

$$W = \begin{pmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, Z = \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad (9)$$

then considering SVD of $E = U \text{diag}(1, 1, 0) V^T$ we can obtain:

$$[\mathbf{t}]_{\times} = U Z U^T, \quad (10)$$

$$R = U W V^T \text{ or } U W^T V^T. \quad (11)$$

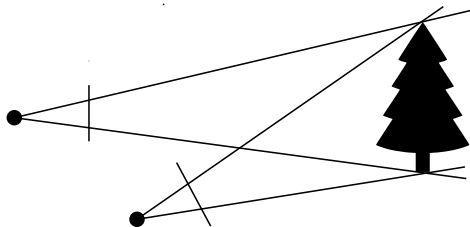
This yields 4 configurations - 2x from choice of R and 2x from changing signs of \mathbf{t} . Only one of these will be valid such that the points are in front both cameras.



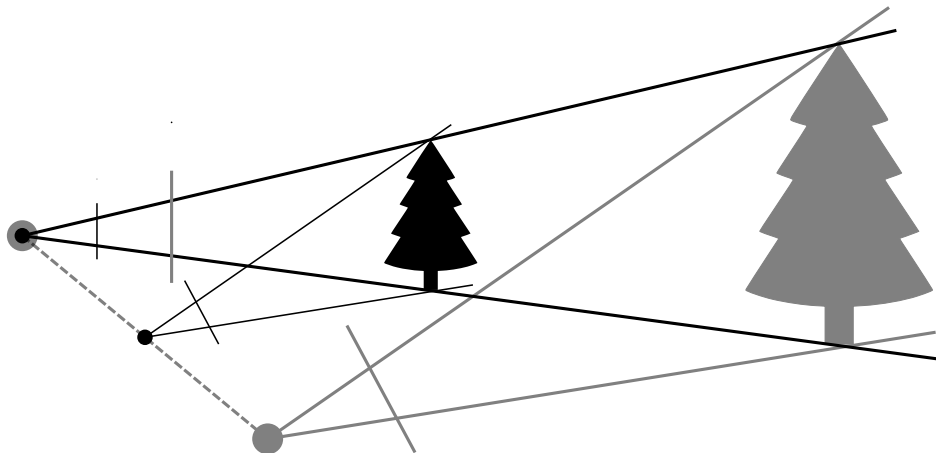
The vector \mathbf{t} we obtained using (10) is always such that $|\mathbf{t}| = 1$. This is due to scale ambiguity.

Suppose that we have a scene with two cameras 10 cm apart imaging an object with size of 5 cm. And the same scene but the cameras are 10 m apart an object with size of 5 m. The images created by the two cameras will be the same. It is therefore impossible to determine the scale of the scene from point correspondences even if we know the camera intrinsics!

Scale Ambiguity



Scale Ambiguity



Why do we care about F ?



You may ask why do we even try to estimate F instead of E in the first place? This is because we can extract additional information about the camera intrinsics from F that are not present in E . For example if you consider that pinhole cameras with zero skew and known principal points then it is possible to obtain the focal lengths of both cameras.

Though in many applications we only calculate E even if the camera intrinsics are not known exactly and we further refine them using non-linear optimization.



We will now discuss algorithms for estimating F . From point correspondences $\mathbf{x}_i \leftrightarrow \mathbf{x}'_i$ we obtain equations:

$$\mathbf{x}'_i{}^T F \mathbf{x}_i = 0. \quad (12)$$

The fundamental matrix also has to satisfy the following:

$$\det(F) = 0, \quad (13)$$

which is a third order polynomial in the entries of F . We will now derive two different approaches: the 7-point minimal algorithm and the 8-point algorithm.



We can rewrite (12) to the form:

$$\mathbf{a}_i^T \mathbf{f} = (x'_i x_i, \quad x' y, \quad x', \quad y' x, \quad y' y, \quad y', \quad x, \quad y, \quad 1) \mathbf{f} = 0, \quad (14)$$

where $\mathbf{x} = (x, y, 1)^T$, $\mathbf{x}' = (x', y', 1)^T$ and \mathbf{f} is the vector representation of F such that:

$$F = \begin{pmatrix} f_1 & f_2 & f_3 \\ f_4 & f_5 & f_6 \\ f_7 & f_8 & f_9 \end{pmatrix}. \quad (15)$$



We can use n equations (14) for n correspondences and rewrite them into:

$$A\mathbf{f} = \begin{pmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{pmatrix} \mathbf{f} = \mathbf{0}. \quad (16)$$

The matrix A has size $n \times 9$ and would have at most rank 8 if we used 8 or more correspondences and they were exact. In that case we could find F via the null-vector of A . However this is not possible in practice.



We will use A generated using 8 or more correspondences. Then we will find the singular vector $\hat{\mathbf{f}}$ corresponding to the smallest eigenvalue using SVD of A . Recall that this minimizes $\|A\hat{\mathbf{f}}\|$ while $\|\hat{\mathbf{f}}\| = 1$.

This does not guarantee that the resulting matrix \hat{F} will be singular. We therefore want to find singular F such that $\|F - \hat{F}\|$ is minimized.



We will use A generated using 8 or more correspondences. Then we will find the singular vector $\hat{\mathbf{f}}$ corresponding to the smallest eigenvalue using SVD of A . Recall that this minimizes $\|A\hat{\mathbf{f}}\|$ while $\|\hat{\mathbf{f}}\| = 1$.

This does not guarantee that the resulting matrix \hat{F} will be singular. We therefore want to find singular F such that $\|F - \hat{F}\|$ is minimized. We can do this using SVD of $\hat{F} = U \text{diag}(\sigma_1, \sigma_2, \sigma_3) V^T$ (we assume $\sigma_1 > \sigma_2 > \sigma_3$). Then the minimizing F is obtained as

$$F = U \text{diag}(\sigma_1, \sigma_2, 0) V^T. \quad (17)$$



The stability of the resulting algorithm depends on the correspondences. It is therefore advisable to perform normalization of the coordinates to improve stability. Instead of using the original correspondences $\mathbf{x}'_i \leftrightarrow \mathbf{x}_i$ we use $\mathbf{q}'_i \leftrightarrow \mathbf{q}_i$ where $\mathbf{q}' = T'\mathbf{x}'$ and $\mathbf{q} = T\mathbf{x}$, where T and T' are transformations performing rotation and translation of the vectors. We then proceed to obtain the fundamental matrix Q . And then we calculate:

$$F = T'^T Q T. \quad (18)$$

T and T' are selected that the centroids of the points are at origin and the root mean square distance of points from origin is $\sqrt{2}$.

7-point Algorithm



We can also find the minimal solution for F given only 7 point correspondences. We can use the same system $A\mathbf{f} = 0$. The matrix A would be rank 7. This means that the null-space is two dimensional. We can use SVD to find two vectors $\mathbf{f}_1, \mathbf{f}_2$ which form a base of the null-space.



We can also find the minimal solution for F given only 7 point correspondences. We can use the same system $A\mathbf{f} = 0$. The matrix A would be rank 7. This means that the null-space is two dimensional. We can use SVD to find two vectors $\mathbf{f}_1, \mathbf{f}_2$ which form a base of the null-space. This means that the final F is a linear combination of F_1 and F_2 . Due to scale ambiguity we can write

$$F = \alpha F + (1 - \alpha)F. \quad (19)$$

To find the final F we solve the rank constraint:

$$\det(F) = \det(\alpha F + (1 - \alpha)F) = 0, \quad (20)$$

which is a cubic equation. This equation may have 1, 3 or no real solutions. In practice we use this algorithm within RANSAC so we can test all of them.



During estimation of F we may encounter some degeneracies. Both algorithms may be unstable when input with configurations close to the degenerate ones.

- Critical surfaces - If the points and both camera centers lie on some types of surfaces (some quadrics) then this might result in degeneracy of the configuration. This could result in multiple solutions for F .
- Points on a plane - If all points lie on a plane then there are multiple solutions to F .
- No translation - If there is no translation then F is not well defined.

In practice we may also have issues when the motion is close to pure translation. In that case it is better to estimate directly for $F = [\mathbf{e}']_{\times}$.



We may use both algorithms to estimate F . In practice we proceed in two steps:

1. We use the 7-point algorithm within RANSAC to find the inlier points.
2. We obtain the final F using the (normalized) 8-point algorithm.

We may also use some iterative optimization methods once we have good initial estimate of F . However, such methods may be slow compared with the 7-point and 8-point algorithms.



We may also want to estimate E . Analogous to 8-point and 7-point algorithms we construct a matrix \hat{A} such that we obtain a system:

$$\hat{A}\mathbf{e} = \mathbf{0}, \quad (21)$$

where \hat{A} is same as A in previous slides but it uses $\hat{\mathbf{x}}$ and $\hat{\mathbf{x}}'$ instead of \mathbf{x} and \mathbf{x}' . We then obtain a four dimensional null-space - $E = \alpha E_1 + \beta E_2 + \gamma E_3 + \delta E_4$. We then solve a system of polynomial equations:

$$\det(E) = 0, \quad (22)$$

$$2EE^TE - \text{tr}(EE^T)E = 0. \quad (23)$$

(23) actually represents 9 equations for the entries of the resulting matrix, but only two of them are algebraically independent. Solving this system is not so straightforward and we will not cover it in this lecture.