

Lecture 14: 2-view Geometry

Lecturer: Luca Carlone

Scribes: -

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor(s).*

This lecture discusses:

- basic principles of 2-view geometry for calibrated cameras (Epipolar constraint, Essential matrix),
- how to compute the essential matrix from pixel correspondences in 2 camera images,
- how to estimate the relative pose (up to scale) between the two cameras from the Essential matrix.

The presentation mostly follows Chapter 5 in [1, Sections 5.1-5.3].

14.1 Epipolar constraint and Essential matrix

From Lecture 12, we know how to compute keypoint correspondences in two images using feature tracking or descriptor-based feature matching. In other words, given a pixel \mathbf{x}_1 in image \mathcal{I}_1 , we are able to compute the corresponding pixel \mathbf{x}_2 in image \mathcal{I}_2 (assuming the two images are picturing the same scene). Note that “corresponding pixels” refers to the fact that the two pixels picture the same 3D point.

In this lecture our goal is to compute the geometry of the two cameras taking images \mathcal{I}_1 and \mathcal{I}_2 (i.e., the relative pose between the cameras) given a number of pixel correspondences. Towards this goal, we take two main assumptions:

- the pixel correspondences are correct, i.e., for every correspondence $(\mathbf{x}_1, \mathbf{x}_2)$, the two pixels *do* represent the same 3D point. We will relax this assumption during the next lecture, since in practice many of the correspondences may be wrong. We also assume that the 3D point does not move.
- the cameras are *calibrated*, i.e., the calibration matrices:

$$\mathbf{K}_1 = \begin{bmatrix} s_{x_1} f_1 & s_{\theta_1} f_1 & o_{x_1} \\ 0 & s_{y_1} f_1 & o_{y_1} \\ 0 & 0 & 1 \end{bmatrix} \quad \mathbf{K}_2 = \begin{bmatrix} s_{x_2} f_2 & s_{\theta_2} f_2 & o_{x_2} \\ 0 & s_{y_2} f_2 & o_{y_2} \\ 0 & 0 & 1 \end{bmatrix} \quad (14.1)$$

of each camera is known. This is an acceptable assumption in robotics, where we can typically calibrate the cameras on our robots and compute \mathbf{K}_1 and \mathbf{K}_2 before deployment (note: for a mobile robot $\mathbf{K}_1 = \mathbf{K}_2$, i.e., both images are collected by the same camera at different time instants).

The perspective projection of the 3D point \mathbf{p}^w to the two cameras (Fig. 14.1) can be written as:

$$p_z^{c_1} \tilde{\mathbf{x}}_1 = \mathbf{K}_1 [\mathbf{R}_w^{c_1} \mathbf{t}_w^{c_1}] \tilde{\mathbf{p}}^w \quad p_z^{c_2} \tilde{\mathbf{x}}_2 = \mathbf{K}_2 [\mathbf{R}_w^{c_2} \mathbf{t}_w^{c_2}] \tilde{\mathbf{p}}^w \quad (14.2)$$

where $(\mathbf{R}_w^{c_i} \mathbf{t}_w^{c_i})$ is the (inverse of) the pose of the camera i in the world frame, $p_z^{c_i}$ is the depth of the point with respect to camera i , and $\tilde{\mathbf{p}}^w$ contains the point coordinates with respect to the world frame (note: homogeneous coordinates).

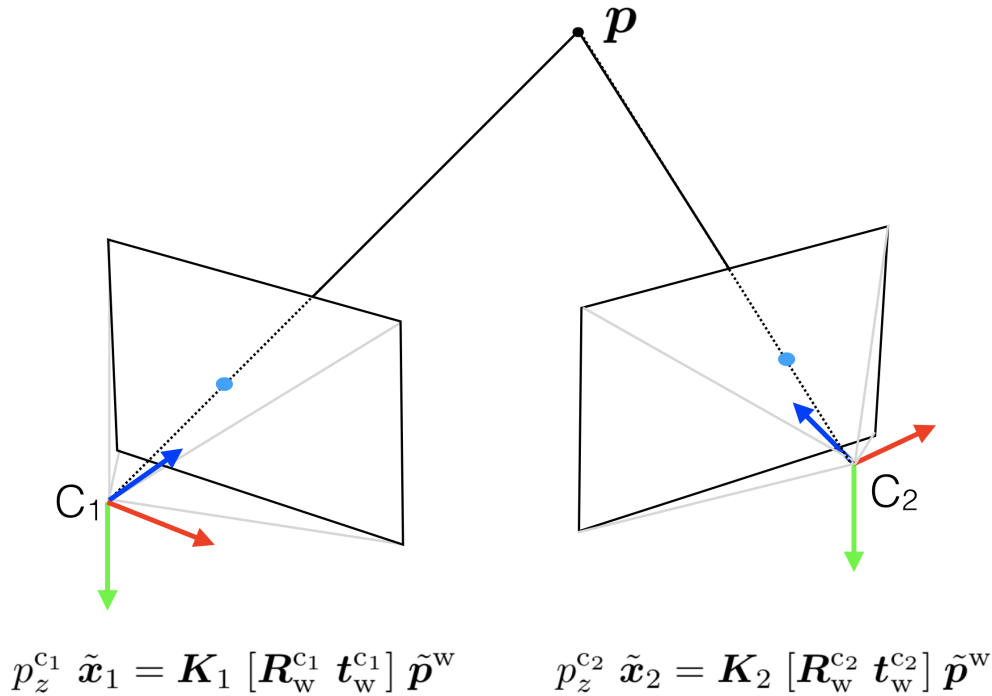


Figure 14.1: 2-view geometry.

Since the point is unknown anyway and we only attempt to compute the relative pose between the two cameras, we may simply assume that c_1 is the w frame, which allows simplifying the expressions above as:

$$d_1 \tilde{x}_1 = \mathbf{K}_1 [\mathbf{I}_3 \mathbf{0}_3] \tilde{\mathbf{p}}^{c_1} = \mathbf{K}_1 \tilde{\mathbf{p}}^{c_1} \quad d_2 \tilde{x}_2 = \mathbf{K}_2 [\mathbf{R}_{c_1}^{c_2} \mathbf{t}_{c_1}^{c_2}] \tilde{\mathbf{p}}^{c_1} \quad (14.3)$$

where to simplify the notation we also defined $d_1 = p_z^{c_1}$ and $d_2 = p_z^{c_2}$.

Since the calibration matrices are known, we can pre-multiply both members of the equation on the left by \mathbf{K}_1^{-1} and both members of the equation on the right by \mathbf{K}_2^{-1} :

$$d_1 \tilde{\mathbf{y}}_1 = \tilde{\mathbf{p}}^{c_1} \quad d_2 \tilde{\mathbf{y}}_2 = \mathbf{R}_{c_1}^{c_2} \tilde{\mathbf{p}}^{c_1} + \mathbf{t}_{c_1}^{c_2} \quad (14.4)$$

where $\tilde{\mathbf{y}}_1 = \mathbf{K}_1^{-1} \tilde{\mathbf{x}}_1$ and $\tilde{\mathbf{y}}_2 = \mathbf{K}_2^{-1} \tilde{\mathbf{x}}_2$ (both still expressed in homogeneous coordinates). $\tilde{\mathbf{y}}_1$ and $\tilde{\mathbf{y}}_2$ are often called “calibrated” pixel coordinates.

Substituting $\tilde{\mathbf{p}}^{c_1}$ from the expression on the left to the right one:

$$d_2 \tilde{\mathbf{y}}_2 = d_1 \mathbf{R}_{c_1}^{c_2} \tilde{\mathbf{y}}_1 + \mathbf{t}_{c_1}^{c_2} \quad (14.5)$$

To simplify notation further, we drop the super- and subscripts for the rotation and translation and write \mathbf{t} (instead of $\mathbf{t}_{c_1}^{c_2}$) and \mathbf{R} (instead of $\mathbf{R}_{c_1}^{c_2}$). This should not cause confusion since these are the only translation and rotation we are going to deal with in this lecture. Therefore, (14.5) becomes:

$$d_2 \tilde{\mathbf{y}}_2 = d_1 \mathbf{R} \tilde{\mathbf{y}}_1 + \mathbf{t} \quad (14.6)$$

Premultiplying both members by $[\mathbf{t}]_{\times}$:

$$d_2 [\mathbf{t}]_{\times} \tilde{\mathbf{y}}_2 = d_1 [\mathbf{t}]_{\times} \mathbf{R} \tilde{\mathbf{y}}_1 \quad (14.7)$$

where we noticed that $[\mathbf{t}]_{\times} \mathbf{t} = \mathbf{0}_3$. Pre-multiplying both members by $\tilde{\mathbf{y}}_2^{\top}$:

$$d_2 \tilde{\mathbf{y}}_2^{\top} [\mathbf{t}]_{\times} \tilde{\mathbf{y}}_2 = d_1 \tilde{\mathbf{y}}_2^{\top} [\mathbf{t}]_{\times} \mathbf{R} \tilde{\mathbf{y}}_1 \quad (14.8)$$

However, $\tilde{\mathbf{y}}_2^{\top} [\mathbf{t}]_{\times} \tilde{\mathbf{y}}_2 = \mathbf{0}_3$ (since $[\mathbf{t}]_{\times} \tilde{\mathbf{y}}_2$ is orthogonal to $\tilde{\mathbf{y}}_2$). Therefore:

$$d_1 \tilde{\mathbf{y}}_2^{\top} [\mathbf{t}]_{\times} \mathbf{R} \tilde{\mathbf{y}}_1 = 0 \quad (14.9)$$

Since d_1 is non-zero, this leads to the *Epipolar constraint*:

$$\boxed{\tilde{\mathbf{y}}_2^{\top} [\mathbf{t}]_{\times} \mathbf{R} \tilde{\mathbf{y}}_1 = 0} \quad (14.10)$$

which relates corresponding pixels in two images. The matrix

$$\boxed{\mathbf{E} = [\mathbf{t}]_{\times} \mathbf{R}} \quad (14.11)$$

is known as the *Essential matrix*.

Definitions:

- epipolar plane: plane passing through the optical centers and the point \mathbf{p}
- epipoles: intersection between the segment connecting the optical centers and the image planes
- epipolar line: line corresponding to the set of pixels $\tilde{\mathbf{y}}_2$ in the second image that satisfy the epipolar constraint for a given pixel $\tilde{\mathbf{y}}_1$ in the first image.

Geometric interpretation. The epipolar constraint can be written as (see slide 12):

$$(\mathbf{t}_{c_1}^{c_2} \times \tilde{\mathbf{y}}_2) \perp (\mathbf{R}_{c_1}^{c_2} \tilde{\mathbf{y}}_1) \iff (\mathbf{t}_{c_1}^{c_2} \times \tilde{\mathbf{y}}_2)^{\top} \mathbf{R}_{c_1}^{c_2} \tilde{\mathbf{y}}_1 = 0 \iff ([\mathbf{t}_{c_1}^{c_2}]_{\times} \tilde{\mathbf{y}}_2)^{\top} \mathbf{R}_{c_1}^{c_2} \tilde{\mathbf{y}}_1 = 0 \iff \tilde{\mathbf{y}}_2^{\top} [\mathbf{t}_{c_1}^{c_2}]_{\times} \mathbf{R}_{c_1}^{c_2} \tilde{\mathbf{y}}_1 = 0$$

Stereo example (slide 14):

$$\tilde{\mathbf{y}}_2^{\top} [\mathbf{t}]_{\times} \mathbf{R} \tilde{\mathbf{y}}_1 = 0 \iff \tilde{\mathbf{y}}_2^{\top} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & b \\ 0 & -b & 0 \end{bmatrix} \tilde{\mathbf{y}}_1 = 0 \iff \tilde{\mathbf{y}}_2^{\top} \begin{bmatrix} 0 \\ b \\ -bv_1 \end{bmatrix} = 0 \iff bv_2 - bv_1 = 0 \iff v_2 = v_1$$

Properties of the essential matrix:

- The epipolar constraint does not constrain the scale of the translation.
- A matrix is an essential matrix *if and only* if it has singular values $\{\sigma, \sigma, 0\}$ (in particular $\sigma = \|\mathbf{t}\|$), see Theorem 5.5 in [1].

Proof. We only prove that the largest eigenvalue of $\mathbf{E} = [\mathbf{t}]_{\times} \mathbf{R}$ is $\lambda_{\max}(\mathbf{E}) = \|\mathbf{t}\|^2$ and the smallest eigenvalue is zero. A complete (but more involved) proof can be found in [1, Thm 5.5].

$$\lambda_{\max}(\mathbf{E}) = \max_{\|\mathbf{d}\|=1} \|\mathbf{E}\mathbf{d}\|^2 = \max_{\|\mathbf{d}\|=1} \mathbf{d}^{\top} \mathbf{E}^{\top} \mathbf{E} \mathbf{d} = \max_{\|\mathbf{d}\|=1} \mathbf{d}^{\top} \mathbf{R}^{\top} [\mathbf{t}]_{\times}^{\top} [\mathbf{t}]_{\times} \mathbf{R} \mathbf{d} \quad (14.12)$$

$$= \max_{\|\mathbf{d}\|=1} \mathbf{d}^{\top} [\mathbf{t}]_{\times}^{\top} [\mathbf{t}]_{\times} \mathbf{d} = \max_{\|\mathbf{d}\|=1} \|[\mathbf{t}]_{\times} \mathbf{d}\|^2 = \max_{\|\mathbf{d}\|=1} \|\mathbf{t} \times \mathbf{d}\|^2 = \|\mathbf{t}\|^2 \quad (14.13)$$

$$\lambda_{\min}(\mathbf{E}) = \min_{\|\mathbf{d}\|=1} \|\mathbf{E}\mathbf{d}\|^2 = \min_{\|\mathbf{d}\|=1} \|\mathbf{t} \times \mathbf{d}\|^2 = 0 \quad (14.14)$$

□

- The space of the essential matrices is called the *Essential space* \mathcal{S}_E (i.e., the space of 3×3 matrices that can be written as $[\mathbf{t}]_{\times} \mathbf{R}$ for some $\mathbf{R} \in \text{SO}(3)$ and $\mathbf{t} \in \mathbb{R}^3$). The projection of a matrix \mathbf{M} onto the Essential space can be computed as prescribed in [1, Thm 5.9]:

$$\arg \min_{\mathbf{E} \in \mathcal{S}_E} \|\mathbf{E} - \mathbf{M}\|_F^2 = \mathbf{U} \begin{bmatrix} \frac{\lambda_1 + \lambda_2}{2} & 0 & 0 \\ 0 & \frac{\lambda_1 + \lambda_2}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{V}^T \quad (14.15)$$

where $\mathbf{M} = \mathbf{U} \text{diag}(\lambda_1, \lambda_2, \lambda_3) \mathbf{V}^T$ is a singular value decomposition of \mathbf{M} .

14.2 How to estimate the Camera Poses from Correspondences?

In this section we address the following problem: given N (calibrated) pixel correspondences, compute the relative pose (up to scale) between the cameras. This is typically done in 2 steps, discussed below:

- use the pixel correspondences and the epipolar constraint to estimate the essential matrix \mathbf{E}
- retrieve the relative pose (\mathbf{R}, \mathbf{t}) between the cameras from the essential matrix \mathbf{E}

14.2.1 Compute the Essential Matrix from Pixel Correspondences

Assume that we are given N (calibrated) pixel correspondences $(\tilde{\mathbf{y}}_{1,k}, \tilde{\mathbf{y}}_{2,k})$ for $k = 1, \dots, N$. As mentioned at the beginning of this document, we assume that there is no outlier (i.e., we do not have wrong correspondences). Each of these correspondences need to satisfy the epipolar constraint (14.10):

$$\tilde{\mathbf{y}}_{2,k}^T \mathbf{E} \tilde{\mathbf{y}}_{1,k} = 0 \quad k = 1, \dots, N \quad (14.16)$$

Noticing that $(\tilde{\mathbf{y}}_{1,k}, \tilde{\mathbf{y}}_{2,k})$ are known pixel values, we realize that these are simply linear equalities. Note that the essential matrix can be only computed up to scale since we can multiply (14.17) by an arbitrary constant without altering the equality. Recalling that $\mathbf{E} = [\mathbf{t}]_{\times} \mathbf{R}$ this means that we can apply an arbitrary scaling to the vector \mathbf{t} without altering the epipolar constraint. This is consistent with the geometric intuition that from a set of camera images we are not able to resolve the scale of the scene, i.e., we are not able to understand if we are observing a small scene from close distance of a large scene from a far distance.

Therefore, it is customary to assume $\|\mathbf{t}\| = 1$, which means we only try to estimate the direction of the translation rather than it's scale.

Eight-point algorithm. In absence of noise, we can compute the essential matrix by solving a linear system. Rearranging the entries of \mathbf{E} in a vector $\mathbf{e} \in \mathbb{R}^9$, the set of linear equations (14.17) can be written as:

$$\mathbf{a}_k^T \mathbf{e} = 0 \quad k = 1, \dots, N \quad (14.17)$$

where \mathbf{a}_k are known vectors whose entries are only function of the pixel correspondences $(\tilde{\mathbf{y}}_{1,k}, \tilde{\mathbf{y}}_{2,k})$.

Stacking the vectors \mathbf{a}_k^T as rows of a matrix \mathbf{A} , the linear equations

$$\mathbf{A} \mathbf{e} = 0 \quad (14.18)$$

For this linear system to admit a unique solution, $\mathbf{A} \in \mathbb{R}^{N \times 9}$ should have rank 8, therefore we need $N = 8$ point correspondences to compute the essential matrix using the linear system (14.18).

By solving the linear system (14.18) and re-arranging the entries of \mathbf{e} into a 3×3 matrix, we obtain the desired essential matrix.

Note that since $\mathbf{A}\mathbf{e} = -\mathbf{A}\mathbf{e} = 0$ both \mathbf{E} and $-\mathbf{E}$ are valid solutions to the linear system (14.18), so we need to consider both as potential essential matrices (we resolve this ambiguity in Section 14.2.2).

Noisy pixel measurements. Since the pixels measurements are typically affected by noise, the solution of the linear system may not be an essential matrix. Therefore, it is common to project the solution onto the essential space using (14.15).

How many correspondences do we really need to estimate \mathbf{E} ? The essential matrix $\mathbf{E} = [\mathbf{t}]_{\times} \mathbf{R}$ is fully defined by a unit vector \mathbf{t} (recall that we imposed $\|\mathbf{t}\| = 1$) and a rotation \mathbf{R} . Therefore, it only has 5 degrees of freedom. Since each correspondence adds a single linear equation, we conclude that we need at least 5 points to estimate the essential matrix.

The eight-point algorithm uses more since it does not leverage the structure of the essential matrix (i.e., we first estimate a generic 3×3 matrix and then we project to the essential space).

Related work indeed provides 7-point, 6-point, and 5-point algorithms. The 5-point algorithm, developed by Nister in [2], is a *minimal* solver.

14.2.2 Retrieve Pose (up to scale) from the Essential Matrix

Theorem 1 (Pose recovery from essential matrix, Thm 5.7 in [1]). *There exist exactly two relative poses (\mathbf{R}, \mathbf{t}) with $\mathbf{R} \in \text{SO}(3)$ and $\mathbf{t} \in \mathbb{R}^3$ corresponding to a nonzero essential matrix \mathbf{E} (i.e., such that $\mathbf{E} = [\mathbf{t}]_{\times} \mathbf{R}$):*

$$\mathbf{t}_1 = \mathbf{U}\mathbf{R}_z(+\pi/2)\boldsymbol{\Sigma}\mathbf{U}^{\top} \quad \mathbf{R}_1 = \mathbf{U}\mathbf{R}_z(+\pi/2)\mathbf{V}^{\top} \quad (14.19)$$

$$\mathbf{t}_2 = \mathbf{U}\mathbf{R}_z(-\pi/2)\boldsymbol{\Sigma}\mathbf{U}^{\top} \quad \mathbf{R}_2 = \mathbf{U}\mathbf{R}_z(-\pi/2)\mathbf{V}^{\top} \quad (14.20)$$

where $\mathbf{E} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{\top}$ is the singular value decomposition of the matrix \mathbf{E} , and $\mathbf{R}_z(+\pi/2)$ is an elementary rotation around the z -axis of an angle $\pi/2$.

Using the matrices \mathbf{E} and $-\mathbf{E}$ we computed in the previous section, we can use Theorem 1 to compute the corresponding pose estimates. Since each essential matrix leads to 2 potential poses, we end up with 4 alternatives. Out of these 4 potential poses, we can find the correct one as follows:

- estimate the position of the 3D points producing the correspondences: all points must satisfy (14.5):

$$d_2\check{\mathbf{y}}_2 = d_1\mathbf{R}\check{\mathbf{y}}_1 + \mathbf{t}$$

from which we can compute d_1 and d_2 . Note that due to the scale ambiguity (we assumed $\|\mathbf{t}\| = 1$) d_1 and d_2 may be a “scaled” version of the true distance from the cameras to the point.

- select the pose for which the reconstructed 3D points are in front of both cameras (the so-called *chirality* constraint).

References

- [1] Y. Ma, S. Soatto, J. Kosecka, and S.S. Sastry. *An Invitation to 3-D Vision*. Springer, 2004.
- [2] D. Nistér. An efficient solution to the five-point relative pose problem. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2003.