

CIS 580, Machine Perception, Spring 2021
Homework 3 Solutions
Due: Friday April 2 2021, 5:59pm

3D Reconstruction from two 2D images

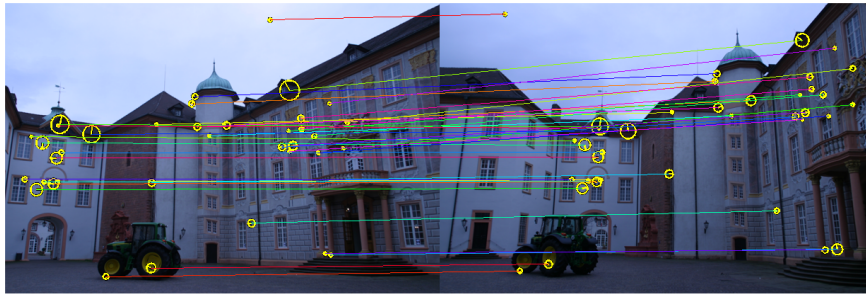


Figure 1: The two images we will use in this problem. A few SIFT matches are shown on top of the images.

Your friend Satsok took some pictures of a castle he visited last summer, and he would love to turn them into a 3D miniature model of the castle. Having heard about your freshly learned computer vision skills, he gives you two pictures (as shown in figure 1) and challenges you to find out where he was standing when he took them and what the 3D scene looks like. He gives you the following details:

- You need to recover the 3D transformation (R, T) between the two views, such that $P_1, P_2 \in R^3$ describe the same scene point in frame 1 and 2, and $P_2 = RP_1 + T$.

We will use the following notations:

- The notation $\mathbf{x} = (x, y, 1)$ (homogeneous metric coordinates) will be used for calibrated projections:

$$P_1 = \lambda_1 \mathbf{x}_1, \quad P_2 = \lambda_2 \mathbf{x}_2$$

- The notation $\mathbf{u} = (u, v, 1)$ (homogeneous pixel coordinates) will be used for uncalibrated projections:

$$\mathbf{u}_1 \sim K\mathbf{x}_1, \quad \mathbf{u}_2 \sim K\mathbf{x}_2$$

$$\mathbf{u} \sim K\mathbf{x} \Leftrightarrow u = fx + u_0, \quad v = fy + v_0$$

- $\text{epi}(\mathbf{x}_1)$ is the epipolar line in camera 2 corresponding to \mathbf{x}_1 .
- He took the two pictures with no zoom and he turned the auto-focus off: the two views share the same

matrix of intrinsics $K = \begin{pmatrix} f & 0 & u_0 \\ 0 & f & v_0 \\ 0 & 0 & 1 \end{pmatrix}$ with $f = 550$ and $(u_0, v_0) = (307.5, 205)$.

- Template code is available [here](#). The images are available [here](#) and [here](#). You will need to submit the completed notebook file.

1.1 Estimation of the essential matrix

You can attempt to reconstruct the scene from the SIFT matchings in the two images (figure 1). The SIFT descriptors are simply computed and matched using OpenCV. Although all the calls to OpenCV for SIFT extraction and matching are already done for you, you might want to take a quick look at OpenCV's tutorials [here](#). Also, quickly read the whole script to understand the overall pipeline.

1. **Least-squares estimation** Complete the function `least_squares_estimation` so that it takes two $N \times 3$ matrices of matched, calibrated points $\mathbf{X}_1, \mathbf{X}_2$ as input and estimates E using the SVD decomposition method as in the 8-point algorithm described in the lecture notes and the “E-matrix” hand-out. Don't forget to project the E you obtain onto the space of essential matrices by redecomposing $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{E})$ and returning $\mathbf{U} \text{diag}(1, 1, 0) \mathbf{V}^T$.
2. **RANSAC estimation** The estimation of E can be made much more robust by selecting sets of points that reach a common agreement: in this section you will implement a basic RANSAC algorithm to eliminate outliers (spurious matchings) and obtain a better estimate of E .
 - (a) One iteration of the algorithm uses the code from the previous question to estimate the essential matrix E based on a random set of 8 correspondences, and then evaluates how good E is for the other pairs $(\mathbf{x}_1, \mathbf{x}_2)$. If E is perfect, for all pairs $(\mathbf{x}_1, \mathbf{x}_2)$, \mathbf{x}_2 lies exactly on the epipolar $\text{epi}(\mathbf{x}_1)$ computed from E and \mathbf{x}_1 . Therefore, we will use the distances to the epipolars as an error measure. Show that the distance of a matching point to the epipolar is the following:

$$d(\mathbf{x}_2, \text{epi}(\mathbf{x}_1))^2 = \frac{(\mathbf{x}_2^T \mathbf{E} \mathbf{x}_1)^2}{\|\widehat{\mathbf{e}}_3 \mathbf{E} \mathbf{x}_1\|^2}.$$

- (b) Complete the function `ransac_estimator` that takes a set of matching points and estimates E using RANSAC. The algorithm steps are the following:
 - Pick a random set of 8 pairs, estimate E using them and compute the individual residuals for all the other pairs $(\mathbf{x}_1, \mathbf{x}_2)$: $d(\mathbf{x}_2, \text{epi}(\mathbf{x}_1))^2 + d(\mathbf{x}_1, \text{epi}(\mathbf{x}_2))^2$.
 - Count how many residuals are lower than $\varepsilon = 10^{-4}$ (consensus set), and if this count is the largest so far, store the current estimate of E as the best estimate so far,
 - Iterate as many times as needed according to the probability of failure
3. **Drawing the epipolar lines** You can now use the essential matrix to plot epipolar lines in the images, as shown in figure 2. Note that E defines epipolars for *calibrated* points $(\mathbf{x}_1, \mathbf{x}_2)$. For the uncalibrated points, we showed in class there is a perfectly identical relationship by substituting the **fundamental matrix** F to E :

$$\mathbf{x}_2^T \mathbf{E} \mathbf{x}_1 = \mathbf{u}_2^T \underbrace{\mathbf{K}^{-T} \mathbf{E} \mathbf{K}^{-1}}_F \mathbf{u}_1 = 0 \Leftrightarrow \mathbf{u}_1 \in \text{epi}_F(\mathbf{u}_2)$$

Complete the function `plot_epipolar_lines` that draws the epipolar lines for a set of matching image points (uncalibrated pixel coordinates). Note that you just need to fill in the equations, the “drawing” part is already done for you.



Figure 2: Some epipolar lines in both images: for a given \mathbf{x}_1 , the matching point \mathbf{x}_2 should be as close as possible to $\text{epi}(\mathbf{x}_1)$ and vice versa.

1.2 Pose recovery and 3D reconstruction

It is now time to recover the transformation R, T between the two cameras. For a given estimate of E , remember that there are two possible solutions for (\widehat{T}, R) (twisted pair ambiguity):

$$\begin{aligned}(\widehat{T}_1, R_1) &= (UR_Z(\frac{\pi}{2})\Sigma U^T, UR_Z(\frac{\pi}{2})^T V^T) \\ (\widehat{T}_2, R_2) &= (UR_Z(-\frac{\pi}{2})\Sigma U^T, UR_Z(-\frac{\pi}{2})^T V^T)\end{aligned}$$

where (U, Σ, V) is the SVD decomposition of E and $R_Z(\frac{\pi}{2}) = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$.

1. Show that the third column of U has associated hat matrix $UR_Z(+\frac{\pi}{2})U'$. Also show that the opposite of the third column of U has associated hat matrix $UR_Z(-\frac{\pi}{2})U'$. Therefore, in the two solutions above, you can output T instead of its hat matrix and simply use the third column of U and its opposite for T_1 and T_2 .
2. Show that the sign switch between $R_Z(+\frac{\pi}{2})$ and $R_Z(-\frac{\pi}{2})$ inside the decompositions of the solutions is equivalent to a left-multiplication of the solutions by $R_T(\pi)$ where T is the third column of U .
3. Complete the function `pose_candidates_from_e` that takes E as an input and returns four possible solutions for (R, T) :
 - the twisted pair for E ,
 - *and* the twisted pair for $-E$. This is because if we find E that satisfies all the epipolar constraints $\mathbf{x}_1^T E \mathbf{x}_2 = 0$ (part 1.1), the opposite matrix $-E$ also verifies the constraint since the right-hand side is zero, and we don't know a priori whether E or $-E$ is the essential matrix we are looking for.
4. **Triangulation** Given a candidate (R_i, T_i) , we can now attempt to reconstruct all the points from the pairs (x_1, x_2) , and check whether (R_i, T_i) is the correct transformation (we need to pick one out of the four candidates): we will pick the candidate (R_i, T_i) that has the highest number of reconstructed

points *in front of both camera*. Consider one of the matchings (x_1, x_2) : we can reconstruct the 3D point by computing the intersection of the two rays coming from camera 1 and 2:

$$\lambda_2 \mathbf{x}_2 = \lambda_1 R \mathbf{x}_1 + T,$$

In practice, the equality doesn't hold (the two rays don't intersect perfectly) and you need to do a least square estimation of λ_1, λ_2 .

For a given pair $(\mathbf{x}_1, \mathbf{x}_2)$, formulate the linear system such that $c\lambda_1\lambda_2$ is its solution. Complete the part of `reconstruct3D` that loops through the four transformation candidates (T, R) and all pairs $(\mathbf{x}_1, \mathbf{x}_2)$, and solves for the depths λ_1, λ_2 .

The selection of the right (T, R) out of the four possibilities is already coded for you but you should take a look at it and make sure you understand it: we simply count for each candidate (T_i, R_i) how many points are in front of both cameras (namely both $\lambda_1, \lambda_2 > 0$) and keep the transformation that achieves the maximum number.

5. You can now run the whole pipeline in the notebook to estimate the transformation (T, R) between the two cameras and the relative positions of the 3D points. Figure 3 shows an example of output from `plot_reconstruction`, shown in the (X, Z) plane (we're on top of the scene, looking down at it like a street map). Note that all depths and the translation (baseline) are up to a scale.
6. Complete the function `show_reprojections` to compare image points in a given camera to the reprojection of images points from the other camera. All you need to do in this function is apply the camera projection model (and be careful about the expression of the 3D transform from 1 to 2 and from 2 to 1).

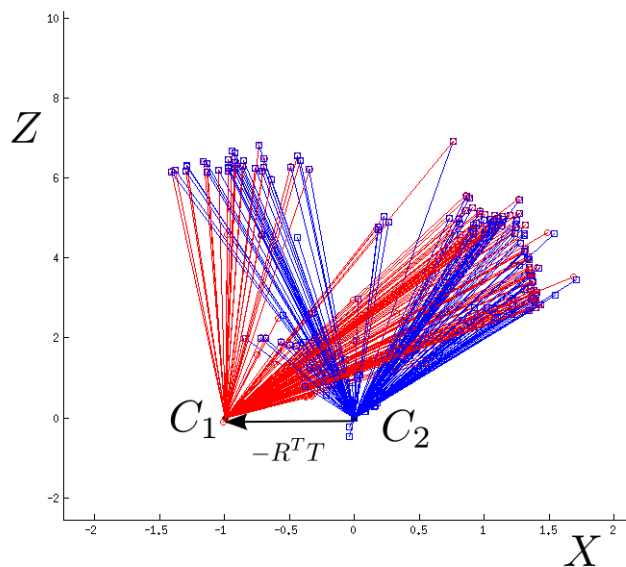


Figure 3: Top view of the the scene, coordinate are expressed in the frame of camera 2 (notice its center is at $(0, 0)$).