

CS280 - Computer Vision - Fall 2013

Homework 2

1 Optical Flow

1. Implement the equations which relate the point wise optical flow to the camera translation and rotation, and depth of points in the scene. Generate visualizations for the following cases (Make sure to look at the function `quiver` in Matlab):
 - (a) Looking forward while driving on the ground plane.
 - (b) Sitting in a train and looking out from a side window.
 - (c) An aircraft heading down towards the ground at an angle of 45° .
 - (d) Camera rotating about the axis $\left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}, 0\right)$.

Deliverables: Your answer in your report should contain - the generated visualizations, code snippet that calculates the optical flow at a point.

2 Error Analysis of Binocular Stereopsis

In this problem we will analyse how error in the measurement of depth varies as a function of depth for a 2 camera system (active like that in Kinect, or passive as in problem 3 below). The way we proceed with estimating the depth Z of a pixel, is by measuring the disparity d for the pixel and then using the formula in Eqn 1.

$$Z = \frac{bf}{d} \quad \text{where, } d = \text{measured disparity, } b = \text{baseline, } f = \text{focal length} \quad (1)$$

Error, δd in the measurement of disparity d , would lead to error δZ in the estimate of Z (you can assume that b and f are known exactly).

1. Derive a formula which relates δZ to Z , δd .
2. `hw2Code/data/kinect.mat` contains the empirical estimate of δZ as a function of Z (and the constants b and f) for a Kinect. Use these to calculate δd for a Kinect.

Deliverables: Your answer in your report should contain - the derivation for part 1, precise and self contained description of how you calculate δd . If you like, you can include a small Matlab code snippet if that makes it any easier for you to explain what you are doing.

3 Multi-view Reconstruction

In this problem, we will see how we can recover the 3D shape of an object given its image from two different viewpoints. In particular, given 2D point correspondences in the 2 images, we will estimate their 3D position in the scene and also estimate the camera positions and orientations.

3.1 Camera Models

Assume $OXYZ$ is the global coordinate system in the 3D scene, while $O_cX_cY_cZ_c$ is the camera coordinate system. Note that, in general these two systems may not be identical. Let $\tilde{\mathbf{X}}$ be a (non-homogeneous) point in the scene (in global coordinates), and $\tilde{\mathbf{X}}_{cam}$ be the coordinates for the same point in the camera coordinate system. It holds that

$$\tilde{\mathbf{X}}_{cam} = R(\tilde{\mathbf{X}} - \tilde{\mathbf{C}}) \quad (2)$$

where $\tilde{\mathbf{C}}$ is the center of the camera in the global coordinate system and R is a 3×3 rotation matrix which describes the orientation of the camera coordinate system with respect to the global coordinate system.

If \mathbf{X} is the 3D point in homogeneous coordinates and \mathbf{x} is the corresponding 2D point in the image, then

$$\mathbf{x} = KR[I | -\tilde{\mathbf{C}}]\mathbf{X} \quad (3)$$

where K is the camera's calibration matrix which is defined by the intrinsic parameters of the camera. We can use the above equation to describe the camera matrix, P : $P = K[R|\mathbf{t}]$, where $\mathbf{t} = -R\tilde{\mathbf{C}}$.

In case of two cameras, it is often easier to assume the global coordinate system to be same as the first camera's coordinate system. In that case, $P_1 = K_1[I|\mathbf{0}]$ and $P_2 = K_2[R|\mathbf{t}]$.

3.2 Fundamental Matrix

The fundamental matrix F is a 3×3 matrix which relates corresponding points in stereo images. Assume $\mathbf{x}_1 \leftrightarrow \mathbf{x}_2$ are corresponding points in an image pair, then the following relation holds:

$$\mathbf{x}_2^t F \mathbf{x}_1 = 0 \quad (4)$$

Note that the fundamental matrix is of rank 2.

3.3 The Essential Matrix

The essential matrix E relates calibrated corresponding image points, and it can be defined from the fundamental matrix as follows:

$$E = K_2^t F K_1 \quad (5)$$

If $P_1 = K_1[I|\mathbf{0}]$ and $P_2 = K_2[R|\mathbf{t}]$, the essential matrix can also be written as

$$E = [\mathbf{t}]_x R \quad (6)$$

where $[\mathbf{t}]_x$ is the representation of the cross product with \mathbf{t} . We can also show that the essential matrix is of rank 2 and has two identical real singular values while the third one is zero (as an optional exercise, you can try to prove it!). For more details refer to Chapter 7 in [1].

Eqn. 6 also tells us that knowing the essential matrix can allow us to estimate the camera orientation and center location as is described below.

3.4 Algorithms

In this section, we describe algorithms that are used to solve various problems related to the task of 3D reconstruction.

3.4.1 Eight-Point Algorithm

This algorithm fits the fundamental matrix defined by corresponding points $\mathbf{x}_1 \leftrightarrow \mathbf{x}_2$ in an image pair. In particular, assume points $\mathbf{x}_1^{(i)} = (x_1^{(i)}, y_1^{(i)})$ in the first image corresponds to points $\mathbf{x}_2^{(i)} = (x_2^{(i)}, y_2^{(i)})$ in the second image for $i = 1, \dots, N$. Note that, $N \geq 8$, we need at least 8 corresponding points.

Normalization. Usually, the point correspondences are not completely accurate. To avoid numerical errors due to noise, we can normalize the points in the image, by translating them by μ (so that the mean of the points is at the origin), and scaling them by σ (so that the mean distance from the origin is a constant (e.g. $\sqrt{2}$)). Assume that the linear transformation described above is represented by the matrices T_1 and T_2 for the two images respectively.

Optimization. Eq. 4 can be rewritten equivalently as $Af = 0$, where f is formed from the entries of F stacked to a 9-vector row-wise and A is a $N \times 9$ dimensional matrix. In particular, the i -th row of A is equal to

$$A_i = [x_1^{(i)} x_2^{(i)} \quad y_1^{(i)} x_2^{(i)} \quad x_2^{(i)} \quad x_1^{(i)} y_2^{(i)} \quad y_1^{(i)} y_2^{(i)} \quad y_2^{(i)} \quad x_1^{(i)} \quad y_1^{(i)} \quad 1] \quad (7)$$

where the point coordinates have already been normalized, i.e. $\mathbf{x} \leftarrow T\mathbf{x}$.

The linear system $Af = 0$ has an exact non-zero solution if $\text{rank}(A) = 8$, which would happen if the point correspondences are exact. But usually due to noise, $\text{rank}(A) > 8$. In that case, there is no exact solution to the linear system and an approximate solution has to be found.

An approximate solution can be found by solving the following optimization problem

$$\min_f \|Af\|_2 \quad \text{s.t.} \quad \|f\|_2 = 1 \quad (8)$$

which can be solved by using the singular value decomposition of matrix A .

The solution F^* found by the approximately solving the problem $Af = 0$, does not guarantee that the fundamental matrix will be of rank 2. We can enforce this constraint, by solving another optimization problem:

$$\min_F \|F - F^*\|_F \quad \text{s.t.} \quad \text{rank}(F) = 2 \quad (9)$$

Again, this problem can be solved using the singular value decomposition of F^* . In particular, if $F^* = USV^t$, where $S = \text{diag}(s_1, s_2, s_3)$ and $s_1 \geq s_2 \geq s_3$ then $F = U\hat{S}V^t$, where $\hat{S} = \text{diag}(s_1, s_2, 0)$, is the matrix that solves the above optimization problem.

Denormalization. After enforcing the rank-2 constraint, we need to remove the normalization such that the fundamental matrix corresponds to points in the actual 2D image space, i.e. $F \leftarrow T_2^t F T_1$.

3.4.2 Estimating Extrinsic Camera Parameters from the Essential Matrix

Another problem is to estimate the extrinsic camera parameters, i.e. R, \mathbf{t} for the second camera, when only the essential matrix E is known.

The operation $[\mathbf{t}]_x$ can also be written as

$$[\mathbf{t}]_x = SZR_{90^\circ}S^t \quad (10)$$

where $S = [\mathbf{s}_0 \quad \mathbf{s}_1 \quad \mathbf{t}]$ is an orthogonal matrix, $Z = \text{diag}(1, 1, 0)$ and R_{90° is the rotation matrix for rotation angle $\theta = 90^\circ$.

From Eq. 6

$$E = [\mathbf{t}]_x R = SZR_{90^\circ}S^t R = U\Sigma V^t \quad (11)$$

Since we know that $\Sigma = \text{diag}(1, 1, 0)$ (remember in homogeneous coordinates we define everything up to a multiplication factor), it holds that $U = S$ and $V = R^t U R_{90^\circ}$.

Thus, the direction of \mathbf{t} as well as R can be determined by the SVD analysis of E . Notice that there is no way we can determine the actual norm of the vector \mathbf{t} from just image point correspondences. Absolute values can only be determined if we have a known reference distance in the scene. In addition, the signs of both \mathbf{t} and R also can not be determined from SVD decomposition. Therefore, this algorithm generates numerous candidate positions and orientations for the camera.

Different combinations of \mathbf{t} and R result in different camera matrix $P_2 = K_2[R|\mathbf{t}]$. After triangulation (3D reconstruction of the point cloud), the combination that gives the largest number of points in front of the image planes is selected.

3.4.3 Triangulation

Triangulation refers to the estimation of the 3D position of a point when the corresponding points in the two image planes are given as well as the parameters of the camera. In particular, assume $\mathbf{x}_1 \leftrightarrow \mathbf{x}_2$ are two corresponding points in the image plane and the camera matrices are P_1 and P_2 respectively. We want to find the 3D point $\mathbf{X} = (X, Y, Z, W)$ such that

$$\mathbf{x}_1 = P_1 \mathbf{X} \quad \text{and} \quad \mathbf{x}_2 = P_2 \mathbf{X} \quad (12)$$

Eq. 12 can also be written as:

$$x_j = \frac{P_{11}^{(j)} X + P_{12}^{(j)} Y + P_{13}^{(j)} Z + P_{14}^{(j)} W}{P_{31}^{(j)} X + P_{32}^{(j)} Y + P_{33}^{(j)} Z + P_{34}^{(j)} W} \quad (13)$$

$$y_j = \frac{P_{21}^{(j)} X + P_{22}^{(j)} Y + P_{23}^{(j)} Z + P_{24}^{(j)} W}{P_{31}^{(j)} X + P_{32}^{(j)} Y + P_{33}^{(j)} Z + P_{34}^{(j)} W} \quad (14)$$

where $j = 1, 2$ are the two camera indices and $\mathbf{x}_j = (x_j, y_j)$.

The above system of non linear equations can be turned into a linear system, if both equations are multiplied by the denominators. The system will be homogeneous if the point in 3D is represented in homogeneous coordinates and thus can be solved via SVD. Alternatively, if we set $W = 1$ then the system is no longer homogenous and can be solved using least squares.

3.5 3D reconstruction

In this assignment, we give you two pairs of images called **house** and **library**. Since the intrinsic parameters of a camera are nowadays known and stored when a picture is taken, we provide the calibration matrices K_1 and K_2 for the two images for both pairs of images. In general, these parameters also need to be estimated via calibration when they are unknown (which is not taught in this class). We also give you corresponding points for both pairs of images. In general, the corresponding points are found by detecting interest points in images and then comparing their descriptors across images to find matches (we will get to this later in the course).

Given a pair of images and their corresponding points as well as the intrinsic parameters for the two cameras, you will have to estimate the 3D positions of the points as well as the camera matrices. In particular, initially you will have to compute the fundamental matrix. You can implement the 8-point algorithm described above or any other algorithm you wish. Subsequently, you will estimate the extrinsic camera parameters of the second camera from the essential matrix. From all the possible combination of parameters, you will keep the one that results in most points in front of the two image planes. Last, you will plot the 3D points as well as the two camera centers. It is up to you how you want to show the point cloud.

3.5.1 Starter Code and Data

`hw2Code/code/reconstruct_3d.m` contains the starter code. This function is the main function and should output the points in the 3D space as well as the camera matrices. The input argument to this function is either `'house'` or `'library'` to select which dataset you want to use. `hw2Code/data` contains the image pairs and camera matrices for the **house** and the **library**.

You will write the following four functions:

1. `fundamental_matrix.m`: This function computes the fundamental matrix F given the data provided. It should return F as well as the residual `res_err`. The residual is defined as the mean squared distance between the points in the two images and the corresponding epipolar lines. Is this what you are directly optimizing using SVD when solving the homogeneous system? If yes, explain. If no, how does the objective relate to the residual? **Deliverables:** Your answer in your report should include - description of how you calculate F , the definition you are using for the residual, how you are calculating the residual, the fundamental matrix and residuals that you get for the 2 sample inputs, and answers to the questions above. Again, if you like, you may include a code snippet if you think it makes your explanation any easier to understand.
2. `find_rotation_translation.m`: This function estimates the extrinsic parameters of the second camera. The function should return a cell array R of all the possible rotation matrices and a cell array t with all the possible translation vectors. **Deliverables:** Your answer in your report should include - all possible choices for t and R for the 2 sample inputs.
3. `find_3d_points.m`: This function reconstructs the 3D point cloud. In particular, it returns a $N \times 3$ array called `points`, where N is the number of the corresponding matches. It also returns the reconstruction error `rec_err`, which is defined as the mean distance between the 2D points and the projected 3D points in the two images. Is this reconstruction error what you are directly optimizing when solving the linear system of equations? If yes, explain. If no, how does the objective of your optimization problem relate to this error? **Deliverables:** Your answer in your report should include - exactly how you calculate the reconstruction error, and reconstruction error for the 2 sample inputs, and answers to the questions above.

4. `plot_3d.m`: This function plots the 3D points in a 3D plot and displays the camera centers for both cameras. Plotting in 3D can be done using the `plot3` command. A plot of the point cloud and the camera centers is enough for the purpose of the assignment. **Deliverables:** Your answer in your report should include - a plot for each of the 2 sample inputs from a reasonable viewpoint.

4 Instructions

1. This HW can be done in groups of 2. We only need 1 submission per group.
2. Please submit through bSpace. In the text box, please put in your name, email address and student ID; and also your partner's if you worked with a partner. Please upload only the following:
 - (a) A single PDF file, containing your answers to all questions (and all supporting images). This is the only thing that we will look at for grading purposes.
3. The HW is due on October 8, 2013, at 11:55 PM.
4. Slip days: As announced already, you have a total of 3 slip days to use through the semester. Slip days can only be used in units of whole days.

References

- [1] Richard Szeliski. *Computer vision: algorithms and applications*. Springer, 2011.