

Mobile Robotics - Final Exam

• Question 1.

1. Questions on Fundamental Matrix [5 points]

- Derive $F e = 0$, where e is the epipole of the 1st camera seen in the second image [2 points]
- If the fundamental matrix between images I_1 and I_2 is F , what is the fundamental matrix between images I_2 and I_1 ? Why are they different? [2 points]
- What is the difference between a fundamental matrix and an essential matrix? How are they related? [1 point]

a) By definition of the fundamental matrix,

$$I_2^T F I_1 = 0$$

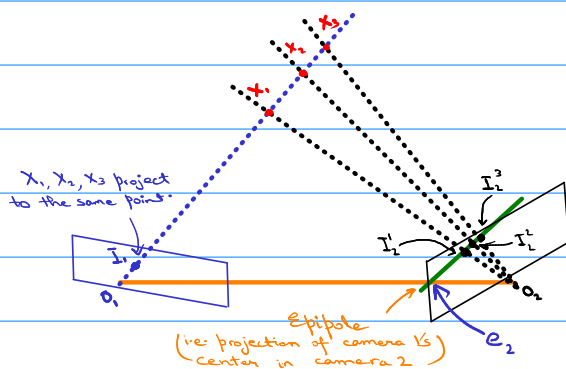
If I_2 and I_1 are images of the same point in the two frames. Such a matrix F is possible since the position of the corresponding point along a line, encapsulated by the nullspace of F , shown below.

For any Image point in 1, its corresponding point in 2 lies on its epipolar line, which also has the epipole. Proof shown below.

$$\Rightarrow e_2^T F I_1 = 0$$

$$\& I_1^T F e_2 = 0$$

$\forall i_1 \leftarrow \text{point in image 1}$
 $\forall i_2 \leftarrow$



I_1, X_1, X_2, X_3, O_1 are coplanar. So, all the points put together, i.e. $I_1, I_2, I_3, X_1, X_2, X_3, O_1, O_2$ are coplanar, this plane is called the Epipolar Plane.

Image plane of the second camera intersects the Epipolar plane, the intersection of 2 planes is a line, called the Epipolar line.

\therefore Any point on the Epipolar plane will lie on this epipolar line.

Given that

$$i_1^T F e = 0 \quad \forall i_1$$

We can conclude that

$$F e = 0$$

b) The definition of the Fundamental Matrix is that:

$$i_1^T F_{12} i_2 = 0 \quad \text{--- ①}$$

for corresponding points i_1 and i_2

Using the same definition for the opposite transform (from 2 to 1), we get:

$$i_2^T F_{21} i_1 = 0 \quad \text{--- ②}$$

Taking transpose of ①

$$i_2^T F_{12}^T i_1 = 0 \quad \text{--- ①}^T$$

Comparing ①^T and ②, we get the equality between the two matrices:

$$F_{21} = F_{12}^T$$

They have to be different because F_{12} takes a full line of image pixels in image 2 to its null space, while F_{21} takes those in image 1.

These projective relations are not symmetric hence the difference.

c) Fundamental Matrix captures epipolar geometry in pixel space, for an uncalibrated camera

Essential Matrix captures the same epipolar geometry in the normalized image space.

This can only be used with a calibrated camera.

Information about camera intrinsics is encapsulated in the Fundamental matrix F , but not in Essential Matrix

Fundamental matrix has 7 degrees of freedom whereas Essential matrix has 5 degrees.

They are related by the following equation (this is by definition and needs no proof)

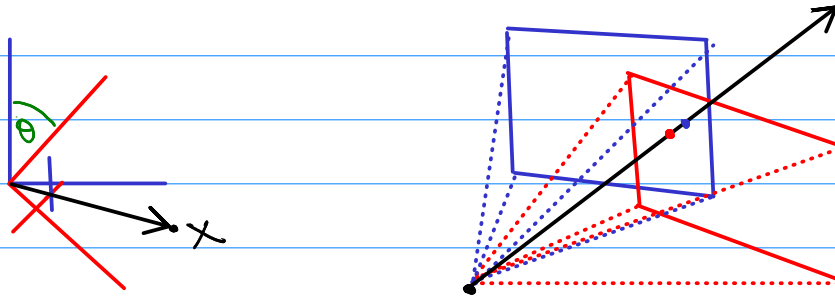
$$\begin{array}{ccccccc} \vec{E} & = & K' & F & K & \longleftarrow & \text{Intrinsic Matrix of Camera 2} \\ \uparrow & & \uparrow & \uparrow & & & \\ \text{Essential Matrix} & & & \text{Fundamental matrix} & & & \\ & & \text{Intrinsic Matrix of Camera 1} & & & & \end{array}$$

• Question 2

2. Questions on relationships between image planes [13 points]

- Derive the homography relation for pure rotation that relates a pixel location x_{i1} in frame I1 to the pixel x_{i2} in frame I2. Does such a homography relation hold when there is a camera translation involved? Explain mathematically your answer [2 + 2 = 4 points]
- What are multiple ways you can relate a pixel x_{i1} in frame I1 to the pixel x_{i2} in frame I2? Be sure to state your assumptions and any additional information used. Write down the equations for those relations [4 points]
- What are the two homographies involved in Stereo Rectification? What does each such homography accomplish? Why would you call them as homographies? Explain with figures and equations [1.5 + 1 + 0.5 + 2 = 5 points]

a)



Given the image point-world point relations of both cameras, we can write:

$$x_B = K_B [I_{3 \times 3} \mid 0_{3 \times 1}] X$$

$$x_R = K_R [R_{3 \times 3} \mid 0_{3 \times 1}] X$$

Rearranging and factoring out R , we get these equivalencies.

$$K_B^{-1} x_B = [I_{3 \times 3} \mid 0_{3 \times 1}] X$$

$$K_R^{-1} x_R = [I_{3 \times 3} \mid 0_{3 \times 1}] R_{3 \times 3} X$$

Equating them both.

$$\Rightarrow R_{3 \times 3} K_B^{-1} x_B = K_R^{-1} x_R$$

$$\Rightarrow x_R = (K_R \quad R_{3 \times 3} \quad K_B^{-1}) x_B$$

→ This is the target homography for the camera rotation

No, the relation does not hold when a camera translation is involved.

This is because our equations become:

$$x_s = K_s [I_{3 \times 3} \mid 0_{3 \times 1}] X$$

$$x_r = K_r [R_{3 \times 3} \mid t_{3 \times 1}] X$$

Now to equate these, we need to invert the projective 3×4 matrix, to recover the 3D point and then equate, only a pseudo-left-inverse will be possible.

So, due to lack of invertibility of $[R_{3 \times 3} \mid t_{3 \times 1}]$ we cannot easily form a homography, and it cannot be written as a simple equation if one exists since we are going to 3-D 4-length vectors.

b) There are 2 classes of relations which relate points in 2 different images

i) Homography: One to One relation of points in one image to those in another.

$$x' = M x$$

point in image 1 Homography point in image 2

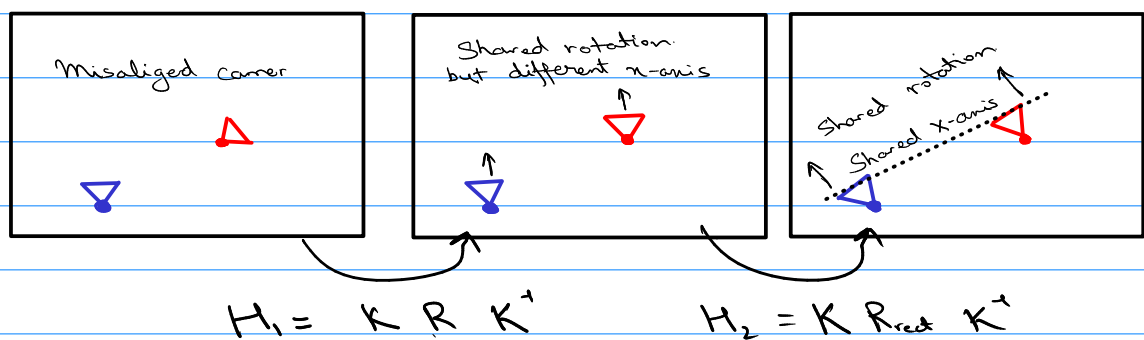
This is possible if:

- i) The two cameras only have a rotation & no translation.
- OR ii) Both images are viewing the same plane from different angles.

- ii) Epipolar geometry: When exact relations don't exist, an image pixel in one image can be mapped to somewhere along a line in the other.
- $x_1^T F x_2 = 0$ as Fundamental matrix
 - $x_1^T E x_2 = 0$ as Essential matrix
- \nearrow point in Image 1 \nwarrow point in Image 2

This holds only if there is non-zero translation between the 2 cameras.

c)



The two homographies are

① Rotational homography $H = K R K^{-1}$

We need both cameras to be in the same rotational frame, but fixed at their initial translation points.

→ Derivation: $x = K [I | 0] X$
 $x' = K [I | 0] R X$
 $\Rightarrow x' = K R K^{-1} x$

② Rotational homography to set epipoles to infinity.

We need to compute the rotation matrix to set epipole E_2 to $[1 \ 0 \ 0]^T$.

The rotation is generated by construction.

$$n' = K R_{\text{rect}} K^{-1} n$$

where $R_{rect} = [r_1 \ r_2 \ r_3]^T$, $\vec{r}_i = \hat{K} e_i = \vec{O}_i - \vec{O}_1$

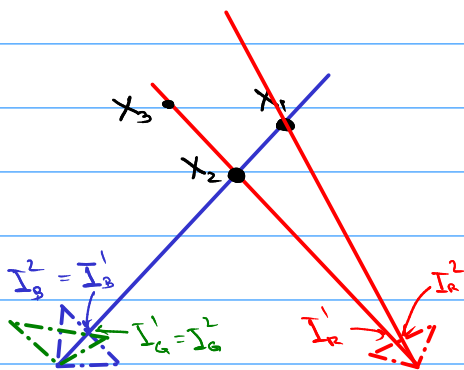
$$\vec{r}_2 = \vec{r}_1 \times [0 \ 0 \ 1]^T \quad \text{and} \quad \vec{r}_3 = \vec{r}_1 \times \vec{r}_2$$

and this takes the epipole e^2 to

$$K^{-1} R_{\text{ref}} K e_2 = [1 \ 0 \ 0]^T$$

c) A Homography is a isomorphism of 2 projective spaces which corresponds to an underlying isomorphism of the vector space that underlies it.

2 points in the real world map to the same point if and only if they are collinear with the projection center.



Any point I_s can be mapped to its equivalent I_a in the rotated camera frame, \therefore

\therefore this is an isomorphism

With translation:

$$\underline{I}_B' \longrightarrow \underline{I}_R^1, \underline{I}_R^2$$

$$I_R' \longrightarrow I_B', I_R''$$

This is not 1-1 so

not an isomorphism.

When space undergoes an isomorphism, i.e. structure & point preserving transform, we call it a homomorphism.

• Question 3

3. Questions on Camera Calibration [4 points]

- a. Write down, in detail, the proof of the DLT algorithm elucidating upon each step. Be sure to highlight the steps in the algorithm that would fail in their purpose if all the correspondences taken lie on a plane. Why is the eigenvector corresponding to the least eigenvalue taken? [2 + 1 + 1 = 4 points]

Direct Linear Transform or DLT helps us estimate the calibration parameters of an uncalibrated camera.

To do so we need to know the 3-D real world coordinates of a few points, and have them identifiable in the image so that we can get the corresponding image coordinates

We know that a projective transform does this mapping

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} P_{11} & P_{12} & P_{13} & P_{14} \\ P_{21} & P_{22} & P_{23} & P_{24} \\ P_{31} & P_{32} & P_{33} & P_{34} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$
$$= \begin{bmatrix} P_{11}X + P_{12}Y + P_{13}Z + P_{14} \\ P_{21}X + P_{22}Y + P_{23}Z + P_{24} \\ P_{31}X + P_{32}Y + P_{33}Z + P_{34} \end{bmatrix}$$

This is a homogeneous equation, we can get 2-equations equating the left and the right:

$$x = \frac{p_{11}x + p_{12}y + p_{13}z + p_{14}}{p_{31}x + p_{32}y + p_{33}z + p_{34}}$$

$$y = \frac{p_{21}x + p_{22}y + p_{23}z + p_{24}}{p_{31}x + p_{32}y + p_{33}z + p_{34}}$$

we can rearrange the terms to write this as a system of linear equations.

$$\begin{bmatrix} -x_i & -y_i & -z_i & -1 & 0 & 0 & 0 & 0 & x_i x_i & x_i y_i & x_i z_i & x_i \\ 0 & 0 & 0 & 0 & -x_i & -y_i & -z_i & -1 & y_i x_i & y_i y_i & y_i z_i & y_i \end{bmatrix} \begin{bmatrix} p_{11} \\ p_{12} \\ p_{13} \\ p_{14} \\ p_{21} \\ \vdots \\ \text{all 12 terms} \end{bmatrix} = 0 \quad \text{--- ①}$$

for all points i

$A \rightarrow$ $\vec{p} \rightarrow$

Since there are 11 parameters in the P vector (12 terms - 1 homogeneity scalar), we need atleast 6 points for this to be fully determined. (2 eqns per point)

However, we have noise in the real world due to which we take an overdetermined set of equations with more than 6 points and try to minimize the residuals.

Writing equation ① again: $\hat{p} = \min_{\vec{p}} A \vec{p}$

We can do SVD on A to get
 $A = U D V^T$,

The last column vector in V is the best solution to $AP=0$, since it's the eigenvector corresponding to the smallest eigenvalue.

This is done because other than the trivial solution $P=0$, all other values of \vec{P} vector can be normalized.

Let $\vec{V}_1, \vec{V}_2, \dots, \vec{V}_n$ be the eigenvectors of A which are the columns of V^T or rows of U .

$$\text{so any } \vec{P} = a_1 \vec{V}_1 + a_2 \vec{V}_2 + \dots + a_n \vec{V}_n$$

$$\Rightarrow A\vec{P} = \lambda_1 a_1 \vec{V}_1 + \lambda_2 a_2 \vec{V}_2 + \dots + \lambda_n a_n \vec{V}_n$$

where $\lambda_1 > \lambda_2 > \dots > \lambda_n$, λ s are eigenvalues

$$\Rightarrow \|A\vec{P}\| \text{ is smallest if } a_1 = a_2 = \dots = a_{n-1} = 0$$

and $a_n \neq 0$, since λ_n is the smallest.

\therefore The smallest eigenvector is used as the solution to \vec{P} , to minimize the norm residuals $\|A\vec{P}\|$.

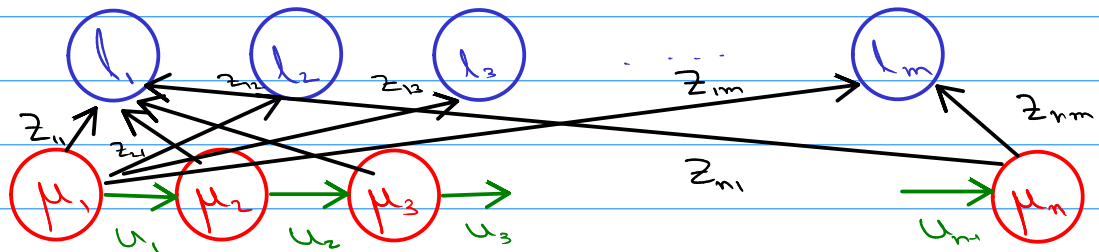
This algorithm will - fail to produce a unique non-0 result because the matrix A would not be full rank. If points are in a plane, we can choose a coordinate frame where $\forall z_i = 0$. $\therefore A$ has at most 3 independent non-0 column vectors, so removing homogeneity P lies somewhere in a 2-D subspace.

Only 1 point in the subspace is valid for the full image, but all points in the subspace correctly map that one plane to the image.

• Question 4

4. **Questions on SLAM/SfM:** For the following questions, clearly write down the variables being solved for and the shape/size of all vectors/jacobians. **[14 points]**
- Given 2D points observed in m observations and the relative pose between any two observations (in the 2D plane - SE(2)), write down the optimization formulation for SAM (Smoothing and Mapping). **[4 points]**
 - You are given a series of RGB images across a trajectory (no additional information) and have to estimate the relative pose between images and map the observed environment. In a systematic and concise way, write down the steps you would take to perform such a monocular SLAM with their mathematical equations. Describe how you obtain your initial estimates and optimise for the trajectory. **[10 points]**

a)



m landmarks (observations per frame)
 n timesteps of motion and observation

$$\mu_i = [x_i \quad y_i \quad \theta_i]$$

We need to be given:

- A motion model eg. $f(\vec{\mu}_i, u_i) \rightarrow \mu_{i+1}$

$$\begin{bmatrix} x_{i+1} \\ y_{i+1} \\ \theta_{i+1} \end{bmatrix} = \begin{bmatrix} x_i \\ y_i \\ \theta_i \end{bmatrix} + \begin{bmatrix} T \cos(\theta_i + \Delta\theta_i) \\ T \sin(\theta_i + \Delta\theta_i) \\ \Delta\theta_i \end{bmatrix}$$
- An observation model eg. $Z_{ij} = \begin{bmatrix} \|\vec{\mu}_i - \vec{l}_j\| \\ \tan^{-1}((x_i - l_{jx}) / (y_i - l_{jy})) \end{bmatrix}$
 $Z_{iA} = h(\vec{\mu}_i, \vec{l}_A)$

Now we can formulate our SAM objective and a linear approximation thereof.

$$L = \sum_{i=0}^n \|f(\bar{p}_i, u_i) - p_{i+1}\|_2^2 + \sum_{l=0}^n \sum_{k=0}^m \|\hat{z}_{ik} - h(l_k, \bar{p}_i)\|_2^2$$

Movement Error
Landmark error

After linearizing $h(l_k, \bar{p}_i) \rightarrow H_{ik} \delta \hat{p}_i + J_k \delta p_{mk} - c_{ik}$

Given this loss term we can compute our residuals Δ \therefore our Jacobian

Jacobian \rightarrow

$$\begin{bmatrix} F_1 & 0 & \dots & 0 & \dots & 0 \\ H_{11} & 0 & \dots & 0 & J_{12} & \dots & 0 \\ H_{12} & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & & & & & \vdots \\ H_{1m} & 0 & \dots & 0 & \dots & J_{1m} & \vdots \\ 0 & F_2 & & & & & \vdots \\ 0 & H_{21} & \dots & 0 & J_{21} & & \vdots \\ 0 & H_{22} & \dots & 0 & 0 & J_{22} & \dots & 0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & H_{nm} & 0 & \dots & J_{nm} \end{bmatrix} \begin{bmatrix} \delta \hat{p}_1 \\ \delta \hat{p}_2 \\ \vdots \\ \delta \hat{p}_m \\ \delta p_{m1} \\ \delta p_{m2} \\ \vdots \\ \delta p_{mn} \end{bmatrix} = \begin{bmatrix} a_1 \\ c_{11} \\ c_{12} \\ \vdots \\ c_{1m} \\ a_2 \\ c_{21} \\ \vdots \\ c_{2m} \\ \vdots \\ c_{nm} \end{bmatrix}$$

Shape of Parameter vector $\rightarrow (3N+2M, 1)$

\therefore N locations of robot with x, y, θ

M landmark locations with x, y

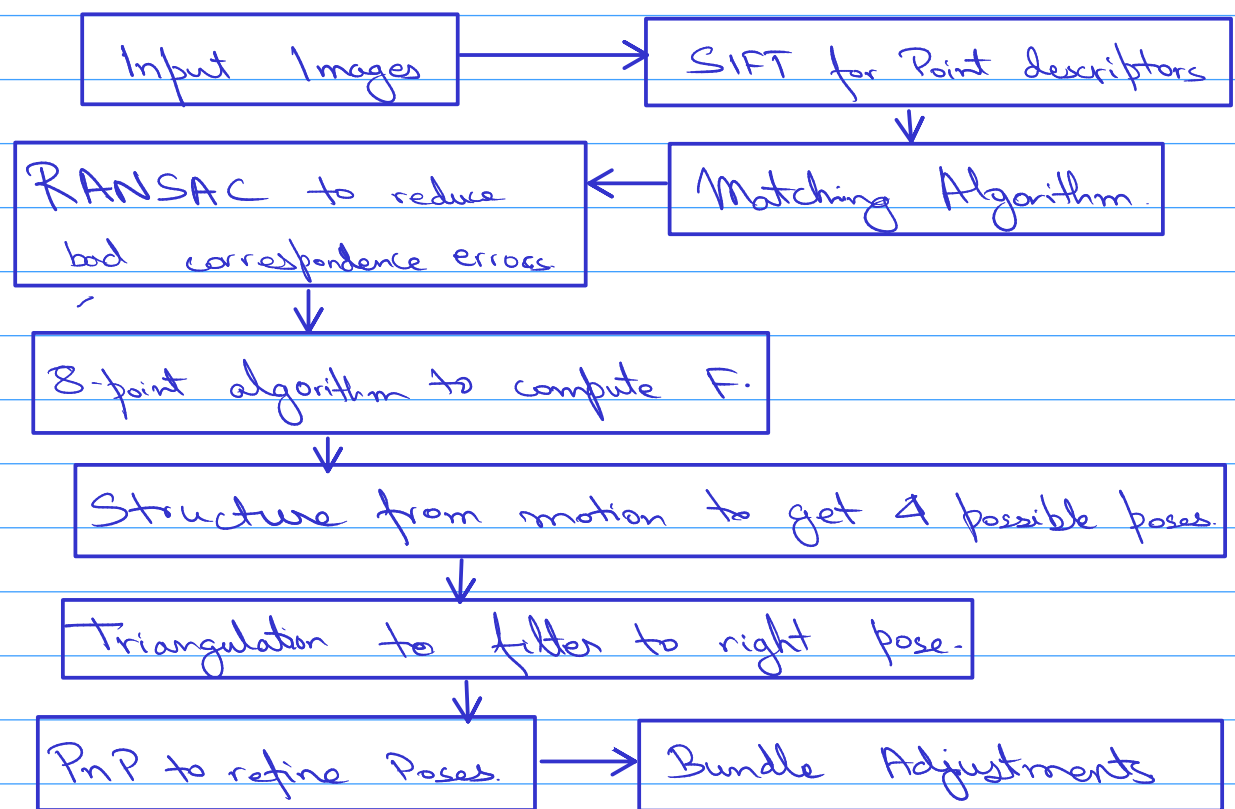
Shape of Residuals vector $\rightarrow (3N+2NM, 1)$

\therefore $(N-1)$ Motion constraints u_i with T and $\Delta \theta$

NM observations of points with even in $T \& \Delta \theta$

\therefore Jacobian has shape $(3N+2M, 2N+2NM)$

b) The steps in the pipeline of computing 3D model from an image trajectory looks as follows.



Now the entire pipeline in more detail.

i) SIFT to generate image descriptors.

Random filters are convolved with the image to generate point descriptors a vector of convolution outputs

$$I * C_i \rightarrow D_i^{N \times M}$$

ii) An all pair point descriptor matching algorithm one example (Rather inefficient) is the Brute Force matcher. we get $f(I_1) = I_2$ where $I_1, I_2 \sim x$ i.e. images of the same point.

iii) Use the 8 point algorithm to compute Fundamental matrix.

We know that for every correspondence x_1, v, x_2
 $x_1^T F x_2 = 0$

we can take F as unknown vector and phrase linear equation

$$A \vec{F} = \vec{0}$$

↓
9 parameter homogeneous matrix

$$A = S_{x_1, x_2} \text{ (Kronecker product)}$$

The smallest eigenvector is an approximation to \vec{F} , so let $U D V^T = A$, we let

$$F_{\text{approx}} = V^T [-1]. \text{ reshape } (3, 3)$$

$\therefore F$ has to be rank deficient,

we compute SVD again $U D V^T = F$.

We drop the lowest eigenvalue to 0 and reconstruct a rank deficient F .

$$\hat{F} = U \begin{bmatrix} D_{11} & 0 & 0 \\ 0 & D_{22} & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T$$

iv) We do RANSAC over correspondences to improve estimate of F .

v) Given our matrix F (or E if we know camera calibration, we use Structure from motion (SfM)'s. We get the following estimates of camera pose from E .

$$\text{Let } E = U D V^T, \quad W = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

So possible configurations are

$$P_1 = [UWV^T \mid U[:,3]]$$

$$P_2 = [UWV^T \mid -U[:,3]]$$

$$P_3 = [UW^T V^T \mid U[:,3]]$$

$$P_4 = [UW^T V^T \mid -U[:,3]]$$

So now we have an ambiguous camera pose.

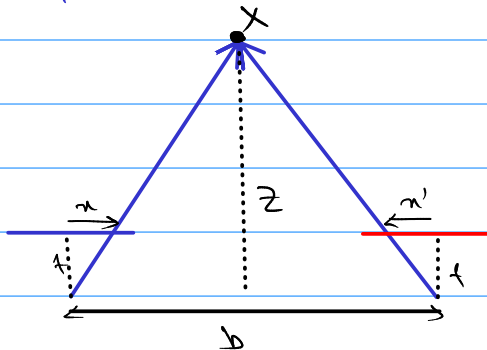
vi) Stereo Rectification is now performed using the two homographies.

$$H = K R K^{-1} \quad \text{on right camera}$$

$$H = K R_{\text{rect}} K^{-1} \quad \text{on both camera.}$$

Now they share scan lines allowing easy triangulation.

vii) Disparity map is generated to get 3D position (i.e. depth information) of points.



$$d = x - x' \\ = \frac{bf}{z}$$

$$\frac{x}{z} = \frac{x}{f} \quad \frac{b-x}{z} = \frac{x'}{f}$$

Using chirality conditions, that camera pose of the 4 poses is selected which has camera seeing the points.

Now we have a point positions as well given the depth information

vii) Perspective from n-points helps us improve our camera pose estimates given the actual 3-D information about those points.

We know correspondences since we generated the 3-D points from 2D.

$$\hat{x}_{3 \times 1} = P_{3 \times 4} X_{4 \times 1}$$

$$\text{Residual} = u - \hat{u}$$

$$\text{Loss} = \|u - PX\|_2^2$$

We optimize the 11 free parameters of P using Levenberg-Marquardt (LM) scheme.

viii) Bundle adjustments: We have the Point estimates and camera pose estimates, we need to refine them.

We are simultaneously refining P_i poses & X_i world map by minimizing Reprojection error

$$\|x_{ij} P_i \bar{X}_j - u_{ij}\|^2$$

by computing Jacobian & using LM algorithm.