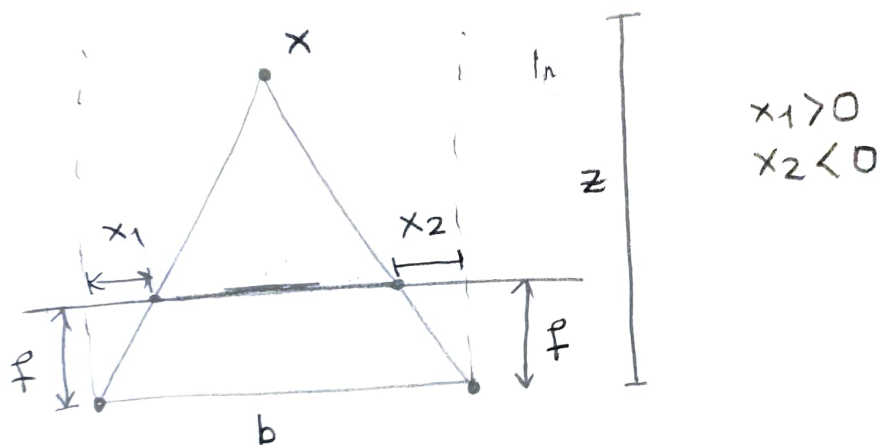## Pen - and - Paper Questions

1. (2 points) What is the difference between depth and disparity? How are they related to each other mathematically?

⟹ Disparity is the horizontal displacement of a point's projections between the left and the right image.

Depth is the horizontal distance to the cameras.



$x_1 > 0$
$x_2 < 0$

In above figure, depth is $z$, and disparity $d = x_1 - x_2$
Using triangles, we can obtain $\dfrac{z-f}{b-d} = \dfrac{z}{b}$

⟹ $zb - bf = zb - zd$

⟹ $zd = bf$  ⟹  $\boxed{d = \dfrac{bf}{z}}$

where $d \to$ disparity
$b \to$ distance between cameras, baseline
$f \to$ focal length
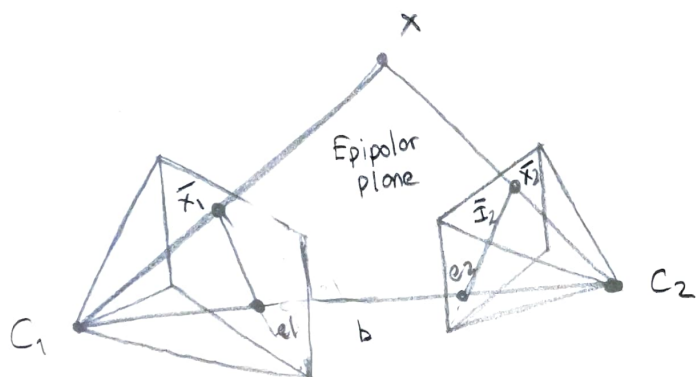$z \to$ depth

2. (2 points) What do we assume to know in calibrated two-view geometry?

We assume we know
- relative position of cameras
- their internal parameters

3. (5 points) Define the following terms related to epipolar geometry.



- Epipolar Line: Epipolar lines shown in above figure with $I_1$ and $I_2$. It's the line connects the projection of 3D point x and the epipole.

- Epipolar Plane: The plane that formed by $C_1, C_2$ (the cameras) and 3D point X.

- Epipole: Epipole is the points that baseline intersects with image planes.

- Projection and Backprojection: Projection is finding the 2D virtual point corresponds to the 3D point. In the figure projecting means finding $\bar{x}_1$ and $\bar{x}_2$ of X. Backprojection is the reverse of it. Finding X from $\bar{x}_1$ and $\bar{x}_2$

- Baseline: In figure, it is b. The line between $C_1$ and $C_2$, the cameras.

4. (4 points) Derive the matrices $M \in SE(3) \subset R^{4\times4}$ representing the following transformations:

- Translation by the vector $T \in R^3$
- Rotation by the rotation matrix $R \in R^{3\times3}$
- Rotation by $R$ followed by the translation $T$
- Translation by $T$ followed by the rotation $R$

Hint: Remember that we can write the transformation matrix $M$ for a given rotation matrix $R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix}$ and a translation vector $T = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}$ as follows: $M = \begin{pmatrix} R & T \\ 0 & 1 \end{pmatrix} = \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix}$

Translation by the vector T

$M = \begin{pmatrix} R & T \\ 0 & 1 \end{pmatrix}$, for just translation by $T$, we can set $R$ to identity matrix $I$.

$\Rightarrow M = \begin{pmatrix} I & T \\ 0 & 1 \end{pmatrix}$ where $I_{3\times3}$ identity matrix.

Rotation by rotation matrix R

$M = \begin{pmatrix} R & T \\ 0 & 1 \end{pmatrix}$, for just rotation by $R$, we can set $T = 0$

$\Rightarrow M = \begin{pmatrix} R & 0 \\ 0 & 1 \end{pmatrix}$

Rotation by R followed by translation T

$\Rightarrow \begin{pmatrix} I & T \\ 0 & 1 \end{pmatrix}\begin{pmatrix} R & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} R & T \\ 0 & 1 \end{pmatrix}$

Translation by T followed by R

$\begin{bmatrix} R & 0 \\ 0 & 1 \end{bmatrix}\begin{bmatrix} I & T \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} R & RT \\ 0 & 1 \end{bmatrix}$

5. (5 points) A classic ambiguity of the perspective projection is that one cannot tell an object from another object that is exactly twice as big but twice as far. Explain why this is true.

Hint: Let $P = (x, y, z)$ be a point on the smaller object and $P' = (x', y', z')$ a point on the larger object. Define $x' = 2x$, $y' = 2y$, $z' = 2z$ and perspective projection as a function $p = \pi(P)$. How does $\pi$ transform the world coordinate $P$ to image coordinate $p$ according to perspective projection? Repeat the same for $P'$ and $p'$

$p = \pi(P)$
_____

$x_s = f \cdot \dfrac{x_c}{z_c}$

$y_s = f \cdot \dfrac{y_c}{z_c}$

$\begin{pmatrix} x \\ y \\ z \end{pmatrix} \rightarrow \begin{pmatrix} f \cdot x/z \\ f \cdot y/z \\ 1 \end{pmatrix}$

$p' = \pi(P')$
_____

$\begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} = \begin{pmatrix} 2x \\ 2y \\ 2z \end{pmatrix} \rightarrow \begin{pmatrix} f \cdot 2x/2z \\ f \cdot 2y/2z \\ 1 \end{pmatrix} = \begin{pmatrix} f \cdot x/z \\ f \cdot y/z \\ 1 \end{pmatrix}$

As it seems from above calculations, the perspective projections of $P$ and $P'$ are same. Therefore the statement is true.

**6. (5 points)**

(a) Given a 3D point $x_w \in \mathbb{R}^3$, such that $x_w = \begin{bmatrix} 3 \\ 2 \\ 6 \end{bmatrix}$, write $x_w$ as an augmented vector $\hat{x}_w$ and as homogenous vector $\tilde{x}_w$

(b) Consider a camera $C$ with pose defined by the rotation and translation matrices, $R = \begin{bmatrix} 0.38 & -0.82 & 0.42 \\ 0.87 & 0.17 & -0.45 \\ 0.29 & 0.54 & 0.78 \end{bmatrix}$ and $t = \begin{bmatrix} 1.3 \\ 2.0 \\ -1.5 \end{bmatrix}$

The camera $C$ has focal lengths $(f_x, f_y) = (785, 786)$, skew $s=0$ and camera center $(c_x, c_y) = (630, 680)$.

Project the point $x_w$ to the image plane of camera $C$.

a/ $\hat{x}_w = \begin{bmatrix} 3 \\ 2 \\ 6 \\ 1 \end{bmatrix}$  we add 1 to make it augmented vector

this augmented vector is homogenous, so

$\tilde{x}_w = \begin{bmatrix} 3 \\ 2 \\ 6 \\ 1 \end{bmatrix}$

b/

$$\tilde{x}_s = [K\ 0]\bar{x}_c = [K\ 0]\begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix}\bar{x}_w$$

$K = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$   $\tilde{x}_s \Rightarrow \begin{bmatrix} 785 & 0 & 680 & 0 \\ 0 & 786 & 680 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}\begin{bmatrix} 0.38 & -0.82 & 0.42 & 1.3 \\ 0.87 & 0.17 & -0.45 & 2.0 \\ 0.29 & 0.54 & 0.78 & -1.5 \\ 0 & 0 & 0 & 1 \end{bmatrix}\begin{bmatrix} 3 \\ 2 \\ 6 \\ 1 \end{bmatrix}$

$\tilde{x}_s = \begin{bmatrix} 6094.6 \\ 5256.9 \\ 5.13 \end{bmatrix}$   we need to divide $z$

$\Rightarrow \begin{bmatrix} 1188 \\ 1024.7 \\ 1 \end{bmatrix}$   so $\begin{pmatrix} 1188 \\ 1024.7 \end{pmatrix}$ is the projection

7) (6 points) We are dealing with a dual-camera setup where the primary camera is positioned to the left of the second camera. The matrices that describe their relative rotation and translation are specified as follows:

- The rotation matrix is identity matrix
- The translation vector is $[110,0]$

a) For $X_1 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$ on image plane of right camera, identify the complete set of corresponding points on the image plane of the left camera, that comply with Epipolar constraint

b) Determine whether the point $X_1 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$ on the image plane of the right camera and the point $X_2 = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$ on image plane of left camera satisfy the Epipolar constraint

c) Given $N$ cameras with $\Pi = \{\Pi_i\}$ their intrinsic and extrinsic parameters, $X_W = \{x_w^P\}$ with $X_w^P \in R^3$ denote a set of $P$ 3D points in world coordinates and $X_S = \{x_i^P\}$ with $X_i^P \in R^2$ denote the image (screen) observations in all $i$ cameras. Give an example of a scenario where the bundle adjustment error is 0?

Hint: In an ideal setting with exact camera parameters and no distortions, observed image points would perfectly align with projections from 3D points, resulting 0 bundle adjustment error.

a) Epipolar constraint $\Rightarrow \tilde{X}_2^T \tilde{E} \tilde{X}_1 = 0$

$\tilde{E} = [t]_x R$

so $\tilde{E} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix}$

$[t]_x = \begin{bmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ -t_y & t_x & 0 \end{bmatrix}$

$\tilde{X}_2^T \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} = 0 \Rightarrow X_2^T \begin{bmatrix} 0 \\ -1 \\ 3 \end{bmatrix} = 0$

$-b + 3c = 0 \quad X_2 = \begin{bmatrix} a \\ 3k \\ k \end{bmatrix} = \begin{bmatrix} a/k \\ 3 \\ 1 \end{bmatrix} \quad \begin{bmatrix} a & b & c \end{bmatrix} \begin{bmatrix} 0 \\ -1 \\ 3 \end{bmatrix} = 0 \rightarrow \begin{bmatrix} a/k \\ 3 \\ 1 \end{bmatrix}$ so $y = 3$

**b/**

we calculated $\tilde{E} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix}$

$\tilde{x_2}^T \tilde{E} \tilde{x_1} = 0$ ? for $x_1 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$ & $x_2 = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$

$\begin{bmatrix} 4 & 4 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix}$

$\Rightarrow \begin{bmatrix} 4 & 4 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \\ 3 \end{bmatrix} = 0 - 4 + 3 = -1 \neq 0$

so does not satisfy epipolar constraint. Also we know from part a, it should be in form $\begin{pmatrix} x \\ 3 \end{pmatrix}$.
y coordinate should be 3, to satisfy epipolar constraint.

**c/**

To get 0 bundle adjustment error, we should know camera intrinsics precisely, there should be no measurement noise, the scene geometry, camera positions, and orientations should be known perfectly without any ambiguity. There should no distortions. Computations should be in infinite precision.

8. (6 points) In a stereo camera setup, we aim to determine the disparity caused by variations in brightness between two cameras set at different exposure levels. We have a sequence of pixel intensity values from a single row in both the left and right camera images as follows:

red

| 19 | 0 | 12 | 2 | 13 | 22 | 20 | 24 | left

| 11 | 19 | 15 | 23 | 18 | 18 | 25 | 24 | right

consider red pixel, whose true disparity is 4. We would like to estimate the disparity of this pixel using block matching with a window size of 3 (1 × 3 as we consider 1D pixel row).

a) Is the true disparity recovered if we use the Sum of squared differences (SSD) similarity metric?

b) Which similarity metric would you use instead considering the brigthness changes between the left & right images? Show that the proposed metric can recover the true disparity for the pixel in a/?

a/ the window spans 22 20 26 on left camera.

disp 5: $\begin{matrix} 22 & 20 & 26 \\ 11 & 19 & 18 \end{matrix}$ → SSD = $(22-11)^2 + (20-19)^2 + (26-18)^2 = 186$

disp 4: $\begin{matrix} 22 & 20 & 26 \\ 19 & 18 & 23 \end{matrix}$ → SSD = $(22-19)^2 + (20-18)^2 + (26-23)^2 = 22$

disp 3: $\begin{matrix} 22 & 20 & 26 \\ 18 & 23 & 18 \end{matrix}$ → SSD = $(22-18)^2 + (20-23)^2 + (26-18)^2 = 89$

disp 2: $\begin{matrix} 22 & 20 & 26 \\ 23 & 18 & 18 \end{matrix}$ → SSD = $(22-23)^2 + (20-18)^2 + (26-18)^2 = 69$

disp 1: $\begin{matrix} 22 & 20 & 26 \\ 18 & 18 & 25 \end{matrix}$ → SSD = $(22-18)^2 + (20-18)^2 + (26-25)^2 = 21$ → smallest disp = 1

disp 0: $\begin{matrix} 22 & 20 & 26 \\ 18 & 25 & 24 \end{matrix}$ → SSD = $(22-18)^2 + (20-25)^2 + (26-24)^2 = 45$

so disparity is 1 estimated did not recovered.

b/ lets use ZNCC as another metric.

22  20  26  → mean $= \frac{22+20+26}{3} = 22.67$

11 19  18 → mean $= 16$  ⇒ ZNCC $= \dfrac{\begin{array}{c}(22-22.67)(11-16)+(20-22.67)(19-16)\\+(26-22.67)(18-16)\end{array}}{\sqrt{0.67^2+2.67^2+3.33^2}\sqrt{5^2+3^2+2^2}}$

$= 0.08$

19 18 23 → m $= 20$  ⇒ ZNCC $= \dfrac{\begin{array}{c}(22-22.67)(19-20)+(20-22.67)(18-20)\\+(26-22.67)(23-20)\end{array}}{\sqrt{0.67^2+2.67^2.333^2}\sqrt{1^2+2^2+3^2}} = \boxed{0.99}$

18 23 18 → m $= 19.7$  ⇒ ZNCC $= \dfrac{(-0.67)(-1.7)+(-2.67)(3.3)+(3.3)(-1.7)}{\sqrt{0.67^2+2.67^2+3.33^2}\sqrt{1.7^2+3.3^2+1.7^2}} = -0.75$

23 18 18 → m $= 19.7$ ⇒ ZNCC $= \dfrac{(-0.67)(3.3)+(-2.67)(-1.7)+(3.3)(-1.7)}{\sqrt{0.67^2+2.67^2+3.33^2}\sqrt{3.3^2+1.7^2+1.7^2}} = -0.19$

18 18 25 → m $= 20.3$  ⇒ ZNCC $= \dfrac{(-0.67)(-2.3)+(-2.67)(-2.3)+(3.3)(4.7)}{\sqrt{0.67^2+2.67^2+3.33^2}\sqrt{2.3^2+2.3^2+4.7^2}} = 0.94$

18 25 24 → m $= 22.3$  ⇒ ZNCC $= \dfrac{(-0.67)(-4.3)+(-2.67)(2.7)+(3.3)(1.7)}{\sqrt{0.67^2+2.67^2+3.33^2}\sqrt{4.3^2+2.7^2+1.7^2}} = 0.06$

so with zncc, estimated disparity is 4.

disparity is recovered.

9. (5 points) Given the focal length $f$ and baseline $b$ for the left camera of a stereo camera setup and the disparity $d$, the depth can be calculated using the simple formula, $z = \frac{fb}{d}$. Show that for an error of $\mathcal{E}_d$ in the estimated disparity, the corresponding error in the obtained depth $\mathcal{E}_z$ is quadratic in depth,

$$\mathcal{E}_z = \frac{z^2}{bf} \mathcal{E}_d.$$

$z = \frac{fb}{d}$      let's say our estimated disparity is

$$\hat{d} = d + \mathcal{E}_d, \text{ so } \hat{z} = \frac{fb}{\hat{d}}.$$

We can define $\mathcal{E}_z$ as $|\hat{z} - z|$, so

$$\mathcal{E}_z = \frac{fb}{d} - \frac{fb}{\hat{d}} = fb\left(\frac{1}{d} - \frac{1}{\hat{d}}\right) = fb\left(\frac{\hat{d} - d}{d\hat{d}}\right)$$

we know $\hat{d} - d = \mathcal{E}_d$, so $\mathcal{E}_z = \frac{fb\,\mathcal{E}_d}{d\hat{d}}$

we know $\hat{d} = \mathcal{E}_d + d$, so $\mathcal{E}_z = \frac{fb\,\mathcal{E}_d}{d(\mathcal{E}_d + d)}$.

Now, we can say $\mathcal{E}_d + d \approx d$ because $\mathcal{E}_d$ is small.

$$\Rightarrow \mathcal{E}_z = \frac{fb\,\mathcal{E}_d}{d^2} = \frac{fb\,\mathcal{E}_d}{(fb)^2/z^2} = \frac{z^2\,\mathcal{E}_d}{fb} \Rightarrow \boxed{\mathcal{E}_z = \frac{z^2}{bf}\mathcal{E}_d}$$

Therefore $\mathcal{E}_z$ is quadratic in depth.