

Computer Vision Notes [1]

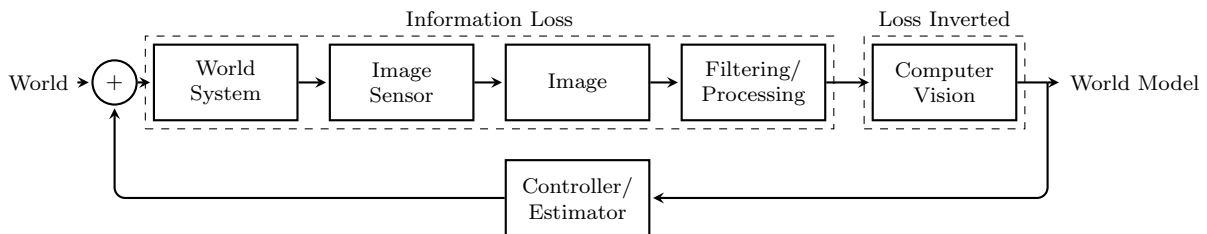
Eric Squires February 16, 2013

Contents

1	Introduction	1
2	Image Formation and Sensing	2
2.1	Image Formation	2
2.1.1	Models of Projection	2
2.1.2	Lenses	2
2.2	Image Sensing	3
2.2.1	Field of View	3
2.2.2	Quantization	3
2.2.3	Transformations	4
2.3	Going from world coordinates to camera pixels: Ψ	6
2.3.1	Deriving the basic form of Ψ	6
2.3.2	Special Topic: Linearizing a Matrix	7
2.3.3	Camera Calibration: Solving for Ψ	9
2.4	Stereo	10
2.4.1	Epipolar Lines	10
2.4.2	Essential vs. Fundamental Matrix	11
2.5	Areas where epipolar lines fail	12

1 Introduction

- Computer Vision - detection of 3d properties (geometric, material) from 2d images (inverse problem)
 - Geometric Properties - size, shape, location
 - Material Properties - radiance, color, texture, material composition
- Computer Vision is not image processing or pattern recognition
 - Pattern Recognition - classifies patterns into a finite number of categories (e.g., is there a person in the picture?)
 - Image processing - producing new image from an old one (often a precursor to computer vision)
- Ways of Approaching Computer Vision: high vs low level, biological vs synthetic
 - Biological - tends to be more complicated/low level
- Block Diagram view
 - Controller - gets something to work (thing in actual world) how we want it to work (model)
 - Estimator - figures out how something (model) works by observing its behavior (actual world)

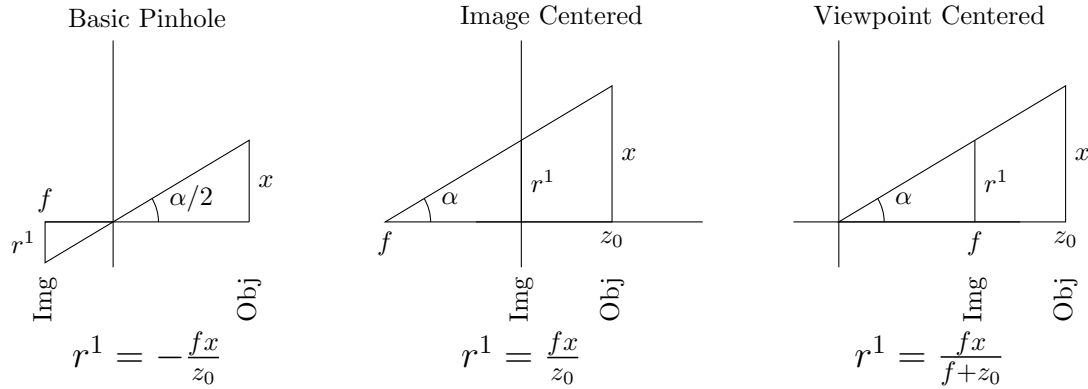


2 Image Formation and Sensing

Steps for image formation and sensing

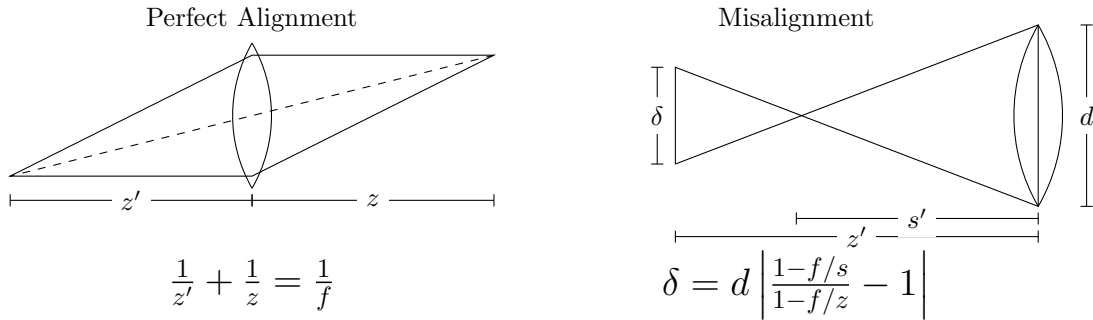
2.1 Image Formation

2.1.1 Models of Projection



- Center of Projection - location of pinhole
- Pinhole problems - does not let in enough light, diffraction (reason for using lenses)
- If $f \gg z$ (e.g., microscope), then $r^1 = \frac{fx}{f+z} \approx x$ (use viewpoint-centered)
- If $f \gg \Delta z$ (e.g., a wall), then $r^1 = \frac{fx}{z+\Delta z} \approx \frac{fx}{z} = mx$ where $m = \frac{f}{z}$

2.1.2 Lenses



To derive the latter equation, define the following:

- z' - actual distance of lens to image plane
- s' - ideal distance of lens to image plane
- z - actual distance to the object
- s - ideal distance to the object

Then by similar triangles, $s'\delta = d(z' - s')$. Then

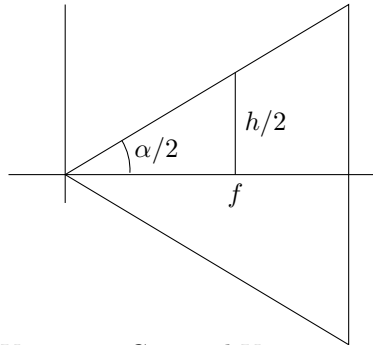
$$\begin{aligned}
 \delta &= \frac{d(z' - s')}{s'} \\
 &= \frac{d|z' - s'|}{s'} && \text{(whether } s' \text{ is too far or too close does not matter)} \\
 &= d \left| \frac{z'}{s'} - 1 \right| \\
 &= d \left| \frac{fz/(z-f)}{fs'/(s-f)} - 1 \right| && \text{(From lens equation: } \frac{1}{z'} + \frac{1}{z} = \frac{1}{f} \Rightarrow z' = \frac{fz}{z-f} \text{ and } s' = \frac{sz}{s-f} \text{)} \\
 &= d \left| \frac{z(s-f)}{s'(z-f)} - 1 \right| \\
 &= d \left| \frac{zs(s-f/s)}{zs(1-f/z)} - 1 \right| \\
 &= d \left| \frac{1-f/s}{1-f/z} - 1 \right|
 \end{aligned}$$

Other interesting (related) points:

- Aperture - a smaller d means less blurring but also less light coming in.
- Depth of Field - range of distances over which objects are focused sufficiently well (e.g., $s \in [z_{min}, z_{max}]$)
- Resolution - higher resolution means lower depth of field (less tolerance for δ)

2.2 Image Sensing

2.2.1 Field of View



Viewpoint Centered View

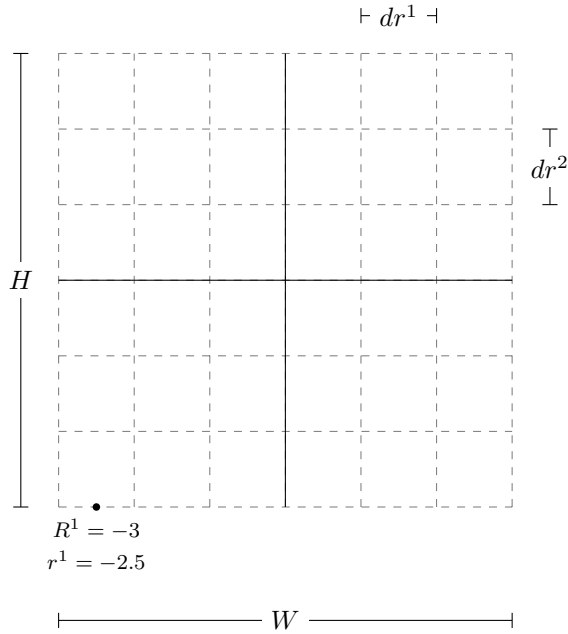
$$\alpha = 2 \arctan \left(\frac{h}{2f} \right) = \text{Field of View}$$

2.2.2 Quantization

Relevant terms¹:

- R^1 - discrete horizontal position (i.e., the pixel #)
- r^1 - continuous horizontal position (i.e., a measurement with a ruler)
- dr^1 - width of a pixel
- W - total image plane width (in pixels)

¹substitute superscript "2" for vertical coordinates



$$R^1 = \lfloor r^1/dr^1 \rfloor + \frac{W}{2} \text{ (} W \text{ even)}$$

$$R^1 = \lfloor r^1/dr^1 \rfloor + \frac{W-1}{2} \text{ (} W \text{ odd)}$$

2.2.3 Transformations

Terminology

1. $T_{BC}^W = \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$ - translation from points B to C in the W (world) reference frame²
2. $R_C^W = \begin{pmatrix} | & | & | \\ \hat{x}_W & \hat{y}_W & \hat{z}_W \\ | & | & | \end{pmatrix}$ - rotation from reference frame W to C (note: $R^T R = I$, $|R| = 1$). To create R_C^W - plot W and C coordinate frames on top of each other. \hat{x}_c is the location of the world's x -axis³ from the perspective of the C frame (see applications below).
3. $g_C^W = \left(\begin{array}{c|c} R_C^W & T_{WC}^W \\ \hline 0 & 1 \end{array} \right)$ - transformation matrix from homogeneous coordinate in reference frame C to same point in W ⁴
4. q_A^W - point A in reference frame W

²i.e., $q_B^W = T_{AB}^W + q_A$

³(x,y,z)

⁴i.e., $q_A^W = \begin{pmatrix} p_A^W \\ 1 \end{pmatrix} = g_C^A q_A^C$

Transformation Summary[2]

Name	Equation	Preserves...	Notes
Translation	$\begin{pmatrix} I & t \\ \mathbf{0}^T & 1 \end{pmatrix}$	angles, lengths, parallel	
Rotation	$\begin{pmatrix} R & \mathbf{0}^T \\ \mathbf{0} & 1 \end{pmatrix}$	angles, lengths, parallel	x, y, z orthonormal
Scaling	$\begin{pmatrix} \alpha I & \mathbf{0}^T \\ \mathbf{0} & 1 \end{pmatrix}$	angles, parallel	
Shear	$\begin{pmatrix} 1 & a & 0 \\ b & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	parallel	
Affine	$\begin{pmatrix} A & \mathbf{0} \\ \mathbf{0}^T & 1 \end{pmatrix}$	parallel	A combines rotation, shear, scale
Projective	$\begin{pmatrix} A & t \\ v & 1 \end{pmatrix}$	straight lines	combines all of the above

3D to 2D Projection [2]: Comparing Orthographic, Para-Perspective, and Perspective Projection

- Orthographic⁵ - removes the z component. Good approximation when $f \gg z$ or $f \gg \Delta z$.

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$

- Para-Perspective - projects on line along line of sight to object center then scales

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{13} \\ a_{13} & a_{13} & a_{13} & a_{13} \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$

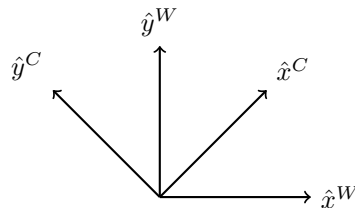
Applications

- Finding R_W^C : Write the coordinates of the world unit x -axis in terms of camera coordinates (same for y , and z). Then

$$R_W^C = \begin{pmatrix} | & | & | \\ \hat{x}_W & \hat{y}_W & \hat{z}_W \\ | & | & | \end{pmatrix}$$

Here is a simple example for finding R_W^C . Given the below picture, we want $\begin{pmatrix} 1 \\ 0 \end{pmatrix}_W \Rightarrow \begin{pmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{pmatrix}_C$

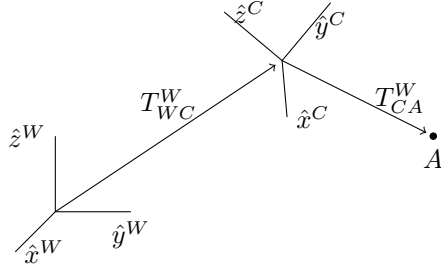
and $\begin{pmatrix} 0 \\ 1 \end{pmatrix}_W \Rightarrow \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}_C$



$$R_W^C = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ -1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

⁵ x and y may also be scaled

- Transform point from perspective of camera to perspective of world



$$q_A^W = R_C^W q_A^C + T_{WC}^W = T_{CA}^W + T_{WC}^W$$

- Turning affine operations into linear ones with homogeneous coordinates

$$\begin{aligned} q_A^W &= \begin{pmatrix} p^W \\ 1 \end{pmatrix} = \begin{pmatrix} R_C^W & \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \\ \hline 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} p^C \\ 1 \end{pmatrix} + \begin{pmatrix} T_{WC}^W \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} R_C^W & \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \\ \hline 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} p^C \\ 1 \end{pmatrix} + \begin{pmatrix} 0 & T_{WC}^W \\ \hline 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} p^C \\ 1 \end{pmatrix} \\ &= \begin{pmatrix} R_C^W & T_{WC}^W \\ \hline 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} p^C \\ 1 \end{pmatrix} = g_C^W q_A^C \end{aligned}$$

- Finding $(g_C^W)^{-1}$.⁶

Inverting the formula from Bullet #1: $q_A^C = (R_C^W)^{-1} q_A^W - (R_C^W)^{-1} T_{WC}^W$

Direct use of formula from Bullet #1: $q_A^C = R_W^C q_A^W + T_{CW}^C$

Combining these two gives: $R_W^C = (R_C^W)^{-1}$ and $T_{CW}^C = -(R_C^W)^{-1} T_{WC}^W$

Using this result as well as the conclusion from bullet #3 gives a new form for $(g_C^W)^{-1}$:

$$q_A^C = g_W^C q_A^W = \begin{pmatrix} R_C^W & T_{WC}^W \\ \hline 0 & 0 & 0 & 1 \end{pmatrix} q_A^W = \begin{pmatrix} (R_C^W)^{-1} & -(R_C^W)^{-1} T_{WC}^W \\ \hline 0 & 0 & 0 & 1 \end{pmatrix} q_A^W = (g_C^W)^{-1} q_A^W$$

2.3 Going from world coordinates to camera pixels: Ψ

2.3.1 Deriving the basic form of Ψ

3 steps

1. Image Sensing ($q^w \rightarrow q^c$): Get points in camera frame (i.e., $q^c = g_w^c q^w$)

$$q^c = \begin{pmatrix} p^c \\ 1 \end{pmatrix} = g_w^c q^w = \begin{pmatrix} R & T \\ \hline 0 & 1 \end{pmatrix} q^w$$

⁶remember that because R is orthonormal, $R^T = R^{-1}$

2. Projection ($q^c \rightarrow \mathbf{r}$ where \mathbf{r} is continuous): Apply perspective projection equations (i.e., $\mathbf{r} = (fx^c/z^c, fy^c/z^c)^T$)⁷

$$\mathbf{r} = \begin{pmatrix} r^1 \\ r^2 \\ 1 \end{pmatrix} = \begin{pmatrix} fx^c/z^c \\ fy^c/z^c \\ 1 \end{pmatrix} \sim \begin{pmatrix} fx^c \\ fy^c \\ z^c \end{pmatrix} = \begin{pmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x^c \\ y^c \\ z^c \end{pmatrix} = \begin{pmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} q^c$$

Now combine this information with that from step 1

$$\mathbf{r} = \begin{pmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \left(\begin{array}{c|c} R & T \\ \hline 0 & 1 \end{array} \right) q^w = \begin{pmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{pmatrix} (R | T) q^w = \Psi_1 (R | T) q^w$$

2 other forms: $\Psi_1 (R|T) q^w = (\Psi_1 R_W^C | \Psi_1 T_W^C) q^w = (\Psi(R_c^w)^T | -\Psi(R_c^w)^T T_c^w) q^w$

3. Quantize Signal - translate,⁸ scale, skew to correct camera abnormalities (e.g., center of focus is not in center of camera, etc). Then round in order to place into buckets

$$\text{Scaling: } \begin{pmatrix} \alpha & 0 & 0 \\ 0 & \alpha & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{Translation: } \begin{pmatrix} 1 & 0 & t_1 \\ 0 & 1 & t_2 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{Skew: } \begin{pmatrix} 1 & \delta & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

When we multiply these equations together, we get a matrix that does all three operations: $\begin{pmatrix} \alpha & \delta & t_1 \\ 0 & \alpha & t_2 \\ 0 & 0 & 1 \end{pmatrix}$.

Multiplying this matrix by what we found in the previous step gives the final Ψ ⁹:

$$\begin{aligned} \begin{pmatrix} \alpha & \delta & t_1 \\ 0 & \alpha & t_2 \\ 0 & 0 & 1 \end{pmatrix} \Psi_1 (R | T) q^w &= \begin{pmatrix} \alpha & \delta & t_1 \\ 0 & \alpha & t_2 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{pmatrix} (R | T) q^w \\ &= \begin{pmatrix} \alpha f_1 & \delta & t_1 \\ 0 & \alpha f_2 & t_2 \\ 0 & 0 & 1 \end{pmatrix} (R | T) q^w = \Psi (R | T) q^w \end{aligned}$$

Note: The rounding step (e.g., floor, ceil) is a non-linear function not captured by Ψ .

2.3.2 Special Topic: Linearizing a Matrix

Suppose you have a stereo camera at $t = t_1$ taking a picture of some object. You then rotate the object and want to know how to transform points on the object at t_1 to the same points at t_2 . Note that the object is rigid so the transformation matrix will be the same for all points. Only from observing points, what is the transformation from t_1 to t_2 ?

⁷note that $(x^c, y^c, z^c)^T = q^c$.

⁸to make the middle the origin of the xy axis, t_1 and t_2 are often set to $W/2$ and $H/2$ respectively. Note that because of projection, translation is just for x and y . There is no z .

⁹note that when the buckets on the camera image sensor are not equal, we get different f_1 and f_2

2.3.3 Camera Calibration: Solving for Ψ

Camera calibration involves identifying camera parameters by taking a picture of a scene where intrinsic calibration solves for Ψ and extrinsic calibration solves for $(R|T)$. In section 2.3.2, a system of the form $q^C = (R|T)q^W$ was solved for $(R|T)$ by linearizing the matrix $(R|T)$. Now multiply both sides by Ψ to get the following:

$$\mathbf{r} = \Psi q^C = \Psi(R|T)q^W$$

Although Ψ was derived in section 2.3.1, people often refer to it in other ways (for reasons listed below):

$$\mathbf{r} = \Psi (R | T) q^w = Dq^w = (\Psi R | \Psi T) q^w = (M | \nu) q^w$$

- Reason for D - In the worst case, $\Psi (R | T)$ has 17 unknowns (5 for Ψ , 9 for R , 3 for T), but using the D form, this comes down to only 12 unknowns. To solve for D (see section 2.3.2), arrange $\mathbf{r}_i = Dq_i^w$ into one of the following forms:

$$0 = \mathbf{r}_i \times Q(q_i)\mathbf{d} \quad 0 = \hat{r}_i Q(q_i)\mathbf{d}$$

where

$$\hat{\mathbf{a}} = \begin{pmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{pmatrix} \quad (1)$$

and $Q(q_i)$ is a matrix after the linearizing process of section 2.3.2. From here, you can solve for D in one of two ways:

- Singular Value Decomposition - given some $m \times n$ matrix A , the **svd** factors A so that

$$A = U\Sigma V^T$$

where U and V^T are orthonormal.¹² Then

$$AV = U\Sigma$$

or looking at just a column at a time gives:

$$A\mathbf{v}_i = \sigma_i \mathbf{u}_i$$

Assuming the problem has a solution, it amounts to finding some \mathbf{v}_i in the null space of A . Fortunately, since we have provided enough points, it is solvable and there should be some $\sigma_k \approx 0$.¹³ Then

$$A\mathbf{v}_k = 0$$

Remember that A stands for $\hat{r}_i Q(q_i)$ and v_k is the \mathbf{d} we are looking for. Putting all this more colloquially, $\mathbf{d} = \mathbf{v}_k$ will always be the right-most vector of V so just run the following code: `[U S V] = svd(Q); V(:,end);`. Then reshape the vector to a 3×4 . Remember at this point though that $\Psi (R | T)$ and Ψ may or may not have a bottom right element of 1. In case it is not 1, divide D by $\sqrt[3]{\det(D)}$.

- Pseudo-Inverse: $\mathbf{d} = \text{pinv}(A)$ or $\mathbf{d} = R \setminus Q$ in MATLAB¹⁴ and do the rescaling in the above bullet.

The downside of using the D matrix is that it needs to be recalibrated every time you move the camera.

¹²orthonormal implies $U^{-1} = U^T$

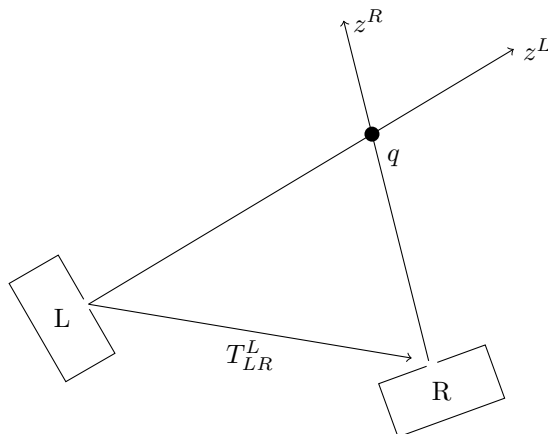
¹³The **svd** will order the singular values from highest (top left) to lowest (bottom right). Also, $k = \min(m, n)$.

¹⁴This method is quite sensitive to floating point arithmetic so it does not work in practice particularly well.

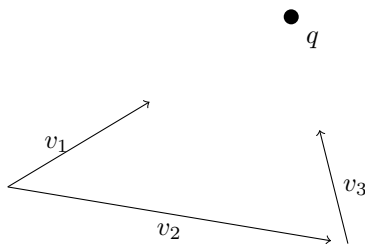
2.4 Stereo

2.4.1 Epipolar Lines

Suppose you have a stereo rig as shown below:



Consider some point projected onto camera L and note that there are an infinite number of points in the world that could have projected to this point. This is z^L in the picture. Similarly, find the same point in camera R and let z^R be the ray representing the set of points that could have projected to this point. Also, assume you know $g_R^L = \left(\begin{array}{c|c} R_R^L & T_R^L \\ \hline 0 & 1 \end{array} \right)$. Note also the following geometric relations:



Because three points determine a plane, we can write

$$0 = (v_3 \times v_2) \cdot v_1$$

and referencing equation (1), we can rewrite this as:

$$\begin{aligned} 0 &= (\hat{v}_3 v_2) \cdot v_1 \\ &= v_1 \cdot \hat{v}_3 v_2 && \text{(dot product is commutative)} \\ &= v_1^T \hat{v}_3 v_2 && \text{(definition of dot product)} \end{aligned}$$

When we compare the two figures in this section, we see an obvious correspondence between z^L , T_{LR}^L ,¹⁵ z^R and v_1 , v_2 , v_3 . There is only one slight change we need to make—we need to rotate z^R into z^L coordinate frame to create z^R . Thus, we are left with:

$$0 = (z^L)^T (\hat{T}_{LR}^L R_R^L) z^R = (z^L)^T E z^R \quad (2)$$

We can now find the set of points z^R that are in the null space of $z^L E$.

¹⁵ $\hat{T} = \begin{pmatrix} 0 & -T^3 & T^2 \\ T^3 & 0 & -T^1 \\ -T^2 & T^1 & 0 \end{pmatrix}$ will be used below

2.4.2 Essential vs. Fundamental Matrix

The Essential matrix is derived in the previous section. To see how the Fundamental Matrix arises, note that $r^L = \Psi z^L$ so that $z^L = \Psi^{-1} r^L$. Then $(z^L)^T = (r^L)^T \Psi^{-T}$.¹⁶ A similar process yields $z^R = \Psi^{-1} r^R$. We can plug these results into equation (2) to find F :

$$0 = (r^L)^T \Psi^{-T} (\hat{T}_{LR}^L R_R^L) \Psi^{-1} r^R = (r^L)^T F r^R$$

This will be a line in the right camera's image. Rather than searching the whole image, we can search for the corresponding point on the line.

To see why this is the case, let $w^T = (r^L)^T F$. Then

$$0 = w^T r^R = w^1 r^{R1} + w^2 r^{R2} + w^3$$

and note that it has the familiar form of a line: $ax + by + c = 0$.

Before contrasting the two matrices, it is helpful to summarize their forms.¹⁷

Matrix	Formula	Image Plane	Ray Relationship
Essential (E)	$\hat{T}R$	$(r^L)^T \Psi^{-T} E \Psi^{-1} r^R = 0$	$(z^L)^T E z^R = 0$
Fundamental (F)	$\Psi^{-T} \hat{T} R \Psi^{-1}$	$(r^L)^T F r^R = 0$	

The difference between the two matrices can be analyzed across four areas: # of parameters, what it maps, solving for the matrix, and using the matrix:

- # Parameters
 - E - 5 parameters (3 for rotation, 2 for translation¹⁸)
 - F - 7 parameters (2 + 2 for epipoles, 3 for homography)
- Solving for the matrix: suppose you have two images taken from stereo cameras of the same object so that there are some corresponding points in both image planes. To solve for F and E , start with their basic formulas and linearize using the method shown in section 2.3.2:¹⁹

$$0 = (r_i^L)^T F r_i^R \quad \Rightarrow \quad 0 = A(r_i^L r_i^R)(\vec{f})$$

$$0 = (r_i^L)^T \Psi^{-T} E \Psi^{-1} r_i^R \quad \Rightarrow \quad 0 = A(r_i^L r_i^R)(\overrightarrow{\psi^{-T} e \psi^{-1}})$$

- E - requires solving for both extrinsic (R, T) and intrinsic parameters (Ψ)
- F - requires solving for only extrinsic (R, T) parameters
- Using the matrix: after reviewing the above table, you will note the following:²⁰
 - E - requires knowing only extrinsic parameters
 - F - requires knowing both intrinsic and extrinsic parameters
- What it maps - the essential (fundamental) matrix maps rays to rays (points to points) as shown in the below figure

¹⁶note that $A^{-T} \equiv (A^T)^{-1}$

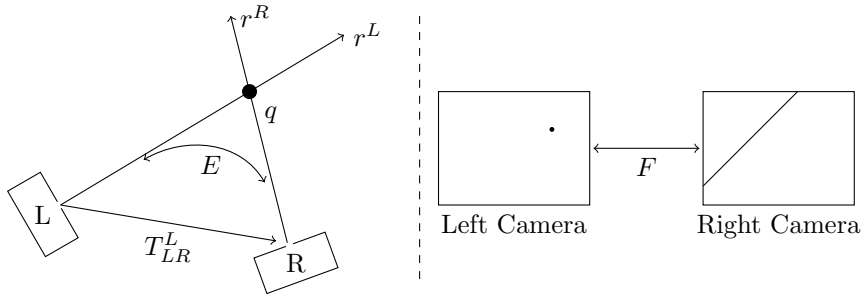
¹⁷It turns out that $E = -R^T \hat{T}$ although the sign generally does not matter because it is in an equation set equal to 0 (equation (2)). To see this, take the transpose of equation (2) so that $0 = 0^T = (r^R)^T (R^T \hat{T}^T r^L$ where $\hat{T} = \begin{pmatrix} 0 & -T^3 & T^2 \\ T^3 & 0 & -T^1 \\ -T^2 & T^1 & 0 \end{pmatrix} = -\hat{T}^T$

(see equation (1)).

¹⁸Magnitude does not matter.

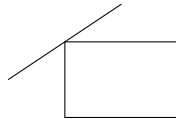
¹⁹The notation \vec{f} means "the linearized version of F "

²⁰Note: if you want to map r^L to r^R , you will always need Ψ . This bullet is mostly just saying the formulas of E and F are different.

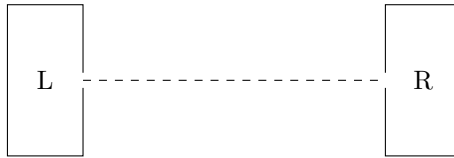


2.5 Areas where epipolar lines fail

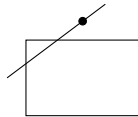
- Epipolar line in one camera does not cross the image plane of the other camera



- Ray from known camera passes through the other camera's optical center



- The point on the epipolar line is outside the field of view



References

- [1] P. Vela, "ECE 4580 class lectures," Spring 2013, (Georgia Institute of Technology).
- [2] R. Szeliski, *Feature detection and matching*. Springer, 2009.