

Stereo!

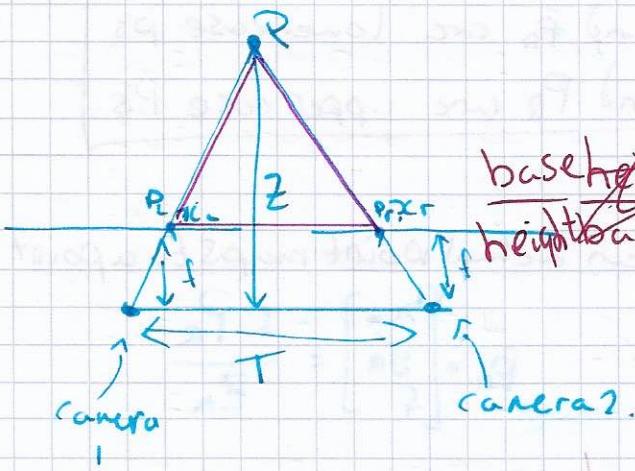
- we can infer the 3D-Structure from 2 or more images taken from different viewpoints.
- The position of each image depends on object in each image depends on its depth.
depth & disparity (positional difference).
The difference between the two images.
- If we have disparity and Viewpoints then we can infer the 3D-Scene structure.

There are 3 key Problems with Stereo.

1. Calibration - we need to determine relative position and orientation of cameras.
2. Correspondence - we need to determine matching points in the Stereo view.
3. Reconstruction - we need to determine the 3D location in a scene of matched points via triangulation.

Triangulation

Simple Stereo:



By Similar Triangles:

$$\frac{\text{base height}}{\text{height base}} = \frac{T}{Z} = \frac{T - x_L + x_R}{Z - f}$$

$$Z = \frac{fT}{x_L - x_R} = \frac{fT}{d}$$

disparity
difference between
location of
Pin left camera (x_L)
and Pin right
camera (x_R).

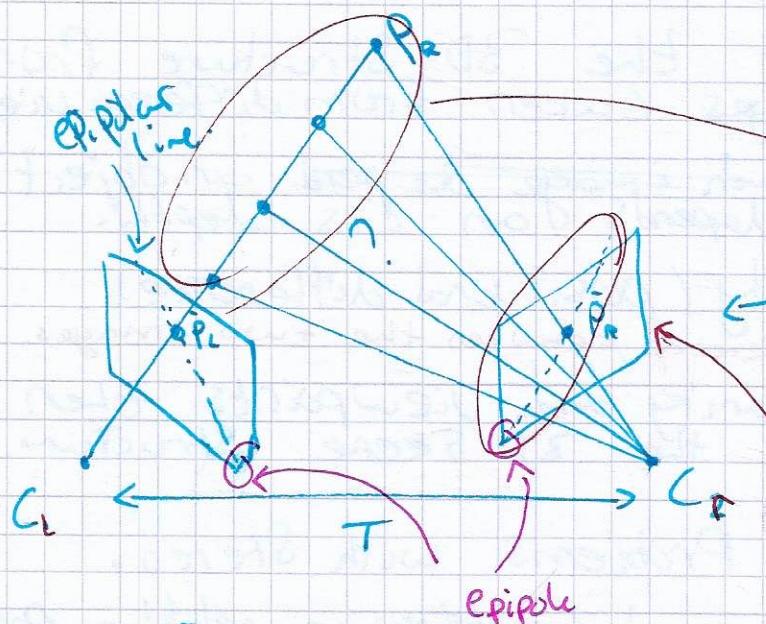
$$\frac{T}{Z} = \frac{T - x_L + x_R}{Z - f}$$

$$TZ - Tf = TZ - Zx_L + Zx_R$$

$$Z(x_L - x_R) = fT$$

$$\frac{Z}{x_L - x_R} = \frac{fT}{Z}$$

Epicolar Geometry.



$$P_L = R^T P_R + T$$

↑ Rotational difference between C_L and C_R

↑ translational difference between C_L and C_R

Epicolar Geometry defines a relationship between two stereo views

- For known viewpoints it constrains matches to lie on along epipolar lines.
- For unknown viewpoints, given the matching viewpoints it is able to estimate the viewpoint.

$$P_L = R^T P_R + T$$

$$P_L - T = R^T P_R$$

$$P_R = R(P_L - T)$$

using previous knowledge of how an actual point maps to a point in the image space.

$$P_{BL} = \begin{bmatrix} x_L \\ y_L \\ z_L \end{bmatrix}$$

$$\bar{P}_L = \begin{bmatrix} x_L \\ y_L \\ f \end{bmatrix} = \frac{f P_L}{z_L}$$

$$\bar{P}_R = \begin{bmatrix} x_R \\ y_R \\ f \end{bmatrix} = \frac{f P_R}{z_R}$$

Epicolar lines

C_L sees the point P but from the image space cannot infer the depth so it could lie anywhere on a line connecting C_L and the actual point P .

C_R sees can view from a different angle. It can 'see' this line. This is the epipolar line.

P from C_R 's perspective must lie on this the epipolar line.

\bar{P}_L and \bar{P}_R are lowercase p's
 P_L and P_R are uppercase P's.

with respect to the left camera

$$\bar{P}_R = P_L - T$$

all vectors P_L, T and $P_L - T$ all lie in the epipolar plane.

$$(P_L - T)^T (T \otimes P_L) = 0 \quad \leftarrow \text{hence } (P_L - T) \text{ and } (T \otimes P_L) \text{ are perpendicular!}$$

↑
in plane

↑
cross product finds a vector perpendicular to plane.

the two vectors T and P_L
since they are both in the plane
wrt Right camera.

the resulting vector
will be perpendicular
to the plane.

$$T \otimes P_L = SP_L$$

$$S = \begin{bmatrix} 0 & -T_x & T_y \\ T_x & 0 & -T_z \\ -T_y & T_z & 0 \end{bmatrix}$$

$$\bar{P}_R^T S P_L = 0$$

$$\bar{P}_R^T R (P_L - T) = 0$$

$$R^T = R^{-1}$$

$$R^T \bar{P}_R = P_L - T$$

$$\bar{P}_R^T R = (P_L - T)^T$$

we call $RS = E$ the Essential matrix

$$\bar{P}_R^T EP_L = 0.$$

$$\bar{P}_L = \frac{f P_L}{Z_L} \quad \bar{P}_R = \frac{f P_R}{Z_R} \Rightarrow \frac{Z_R}{f} \bar{P}_R^T E \frac{Z_L}{f} \bar{P}_L = 0$$

$$\bar{P}_R^T E \bar{P}_L = 0.$$

now let's

$$\bar{u}_L = E \bar{p}_L = \begin{bmatrix} u_{L1} \\ u_{L2} \\ u_{L3} \end{bmatrix}$$

then,

$$\bar{P}_R^T E \bar{P}_L = \bar{P}_R^T \bar{u}_L = x_R u_{L1} + y_R u_{L2} + f u_{L3} = 0.$$

This is the equation of the epipolar line in the right image.

We can also do the same in similar method to get the equation of the epipolar line in the right image. / left image.

Pixel values represent light intensity within a small region of the image plane of size $S_x \times S_y$.

Each pixel has coordinates (\hat{x}_c, \hat{y})

The image coordinates are $x = S_{xc}(\hat{x}_c - \bar{\alpha}_x)$

$$y = S_y(\hat{y} - \bar{\alpha}_y)$$

$$\text{so } \hat{P}_L = \begin{bmatrix} x_c \\ y_c \\ f \end{bmatrix} = \hat{M}_L \begin{bmatrix} \hat{x}_c \\ \hat{y} \\ f \end{bmatrix} = M_L \hat{P}_L$$

we also have

$$\hat{P}_R^T E \hat{P}_L = 0$$

$$\hat{P}_R^T E M_L \hat{P}_L = 0$$

$$\hat{P}_R^T M_R^T E M_L \hat{P}_L = 0$$

$$\hat{P}_R^T F \hat{P}_L = 0$$

$$F = M_R^T E M_L$$

F is the fundamental matrix

We can also estimate F if we are given a set of point correspondences

$$\text{ie } \hat{P}_R^T(:, i) F \hat{P}_L(:, i) = 0 \quad \text{for } i=1\dots N$$

This

$$Af = 0$$

components of f .

$N \times 9$ matrix defined by correspondence vectors $\hat{P}_R(:, i), \hat{P}_L(:, i)$.

We can solve for f using Single Value Decomposition?

3D-Reconstruction.

Where does this come from?

$$\text{Find } a, b, c \text{ st. } a\bar{p}_L - bR^T p_R - T - c(p_L \otimes R^T p_R) = 0.$$

Given corresponding points we know p_L, p_R .
Given calibrated views we know, R, T .

$$a \begin{bmatrix} p_L \\ 3 \times 1 \end{bmatrix} - b \begin{bmatrix} R^T p_R \\ 3 \times 1 \end{bmatrix} - c \begin{bmatrix} p_L \otimes R^T p_R \\ 3 \times 1 \end{bmatrix} = T$$

H is a
 3×3 matrix

$$H = \begin{bmatrix} \vdots & \vdots & \vdots \\ p_L & R^T p_R & p_L \otimes R^T p_R \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} = T \Rightarrow \begin{bmatrix} a \\ b \\ c \end{bmatrix} = H^{-1} T.$$

$$\hat{p} = (a p_L + b R^T p_R + T) / 2. ?$$

Stereo Correspondence

Right camera perspective with R^T and T being added to map the point to the same space as the left camera.

- How do we match up points from each camera?

- What do we match up between cameras?

Pixels, Regions, Feature Points, edges and boundaries, Objects.

Pixels - dense but not distinctive ← Could have lots of similar pixels.

Regions - dense but more distinctive

Feature Points - Sparse but distinctive

Edges - Sparse but ambiguous Lots of few edges, but edges are hard to match up

Objects - Sparse, distinctive but hard to identify and match. ↑

(Could detect a flower in a flowerbed) but how sure are you it's the same flower.

Where to Search?

- uncalibrated views: correct matches could lie anywhere in other views.
- known fundamental matrix - matches constrained to lie along the epipolar lines.

How to Compare Elements?

This depends on what you are trying to match.

You could use:

- Pixel intensities (~~RGB~~)
- Local colour variation.
- Local Structure Variation.
- Position and angle of edges
- Shape and pose.

There are several challenges when comparing edges.

- lighting and noise effects.
- Perspective transformations.
- Occlusions. - things may cover an object from another perspective.
- Repetitive Structure. (chess board) how do you know which area we are in?.

Region Based Methods

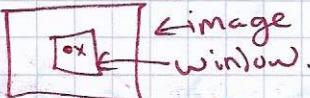
Compare pixel values in image regions

- For a region in the left image compute similarity with regions of same size in the right image.
- The corresponding region is the centre of the most similar region.

Stereo image pair $I_L(x, y)$ and $I_R(x, y)$.

For each pixel find the disparity $d = (d_x, d_y)$ which minimises/maxes the cost function.
(find the offset which maximises the similarity between pixels).

$$C(x, y, d) = \sum_{(i, j) \in W(x, y)} \text{Sim}(I_L(x+i, y+j), I_R(x+i+d_x, y+j+d_y))$$

$W(x, y)$ is the window of pixels around (x, y) . 

There are many similarity measures we could use:

- Sum of squared differences.

$$\text{Sim}(u, v) = (u - v)^2$$

- Similar pixel count

$$\text{Sim}(u, v) = \begin{cases} 1 & \text{if } |u-v| < T \\ 0 & \text{else.} \end{cases}$$

threshold
similarity
difference.

- Normalised cross correlation:

$$\text{Sim}(u, v) = (u - \bar{u})(v - \bar{v}) / N_u N_v$$

$$\bar{u} = \frac{1}{|W(x, y)|} \sum_{k, l \in W(x, y)} u(k, l)$$

mean.

$$N_u = \sqrt{\sum_{k, l \in W(x, y)} (u - \bar{u})^2} \text{ of variance.}$$

Feature Based Methods

- This restricts the search to a much sparser set of features.
- we can find distinct points in each view and then match these points by comparing pixels or image descriptors in local regions about each point.

To example Harris corner detector
RANSAC for uncalibrated matching, matching using the Scaled-invariant feature transform.

Harris Corner detector

This uses the covariance of spatial gradient vectors within the region W .

$$A = \sum_{x,y \in W} \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}$$

Saliency = how important

The eigen values λ_1 and λ_2 of A indicate the spread of gradients in the region. For example 2 high values indicate a busy region.

The Saliency metric $\text{Sal} = \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2}$

If both λ_1 and λ_2 are large then Sal is large

If λ_1 or λ_2 is small then Sal is small.
Can threshold Saliency to get key points.

Calibrated Feature Matching

- We determine the Salient Points in each image.

For each point in one image find the best matching point in the other image.

We can determine the best matching point by using region based matching ~~along the~~ around the Salient points. We also only search around the ~~sets~~ epipolar lines.

Because they tell us where the point should lie based on the other camera's perspective.

Uncalibrated Matching

- When we have calibrated cameras we are able to use the epipolar lines
- However when we do not have calibrated cameras we cannot use the epipolar lines, so we need another approach.

One such approach is to use RANSAC - Random Sample Consensus.

1. Find potential correspondences

2. Select a subset of matches at random (minimum 7).

3. Compute F from subset.

4. Assess how good F is for all other correspondences.

5. Repeat until you find the best F.

Basically k-fold cross validation. k - Number of ~~times an~~ times an F is calculated.

Scale-Invariant Feature Transform

There are 2 main components: key

1. Scale invariant detection of salient points.

2. Match using highly distinct local descriptors.

1. key (Salient) point detection.

- extreme max/min in the difference of Gaussian blurred versions of image

Dog tree.

↑ Difference of Gaussians

- Points imaged at different resolutions appear at different levels of Dog tree.

Scale invariance

Generate images with different levels of blur. calculate pairwise differences

reduces our search space.

↑
Or we will have to do a lot of comparisons.
BAD!

2. Spatial Gradient Descriptors:

- Built from histograms of spatial gradients in local neighbourhood.
- Good rotation & and perspective wrap properties.