

Introduction to 3D image processing

Topics:

1. Camera model and 3D geometry
2. Depth estimation
3. Point cloud
4. SLAM (simultaneous localization and representation model)
5. NERF model

Overview of 3D image processing

Main object: RGB-D Images

- RGB - 3 color channels
- D - depth, distance from the source to the object in a pixel

Image with depth also called depths map

Auxiliary object: point clouds

Unordered set of 3D space points. Every point associated with:

- XYZ coordinates
- Optionally: color, class

Represents a single object in a scene

Point clouds tasks

- Classification
- Part segmentation
- Completion

- Semantic segmentation

Depth estimation

Input

- Monocular image / video

- Stereo pairs

- Sparse depth map

Output

- Dense depth map

Simultaneous localization and mapping

Approaches:

- 3D modeling

- Neural scene fitting

Approaches:

- 3D modeling

- Neural scene fitting

New View Generation

Task: for a given single view or multiple views with known camera parameters to generate a new view defined by camera's position and direction

Approaches:

- 3D modeling

- Neural scene fitting

Sources of 3D data

① Hardware approach: RGB-D sensors

Popular depth sensors

1. Radar
2. Depth camera
3. LIDAR

Radar uses radio waves to determine distances to various objects.

Depending on frequency the range changes.

Advantages:

- + Relatively robust to illumination, weather
- + Cheap

Disadvantages:

- Sparse
- No visual information
- Hard to properly detect small or stationary objects

Depth camera

A setup which estimates depths

Advantages:

- + Sees colors due to RGB camera involved

Disadvantages:

- Sparse depth
- Some cameras have small range

- Heavily affected by the illumination and its intensity

LIDAR

Measures the amount of time of travel to the reflected light from a laser

Advantages:

- + High accuracy for low distance objects

Disadvantages:

- High cost
- "Sees" shapes, not colors

- Affected by weather

② Use computer vision to obtain RGB-D image

- Stereo image

- Single image

- Video

③ Computer vision with hardware

- Turning sparse map into dense map
- Improving map resolution
- "Filling" holes
- Reducing noise
- Combine color and depth info to improve analysis

- Basics of 3D image processing
- Camera model
- Perception

An object which are sensor. Each generates

different angles

by rays from multiple points of the source object resulting in blurry image.

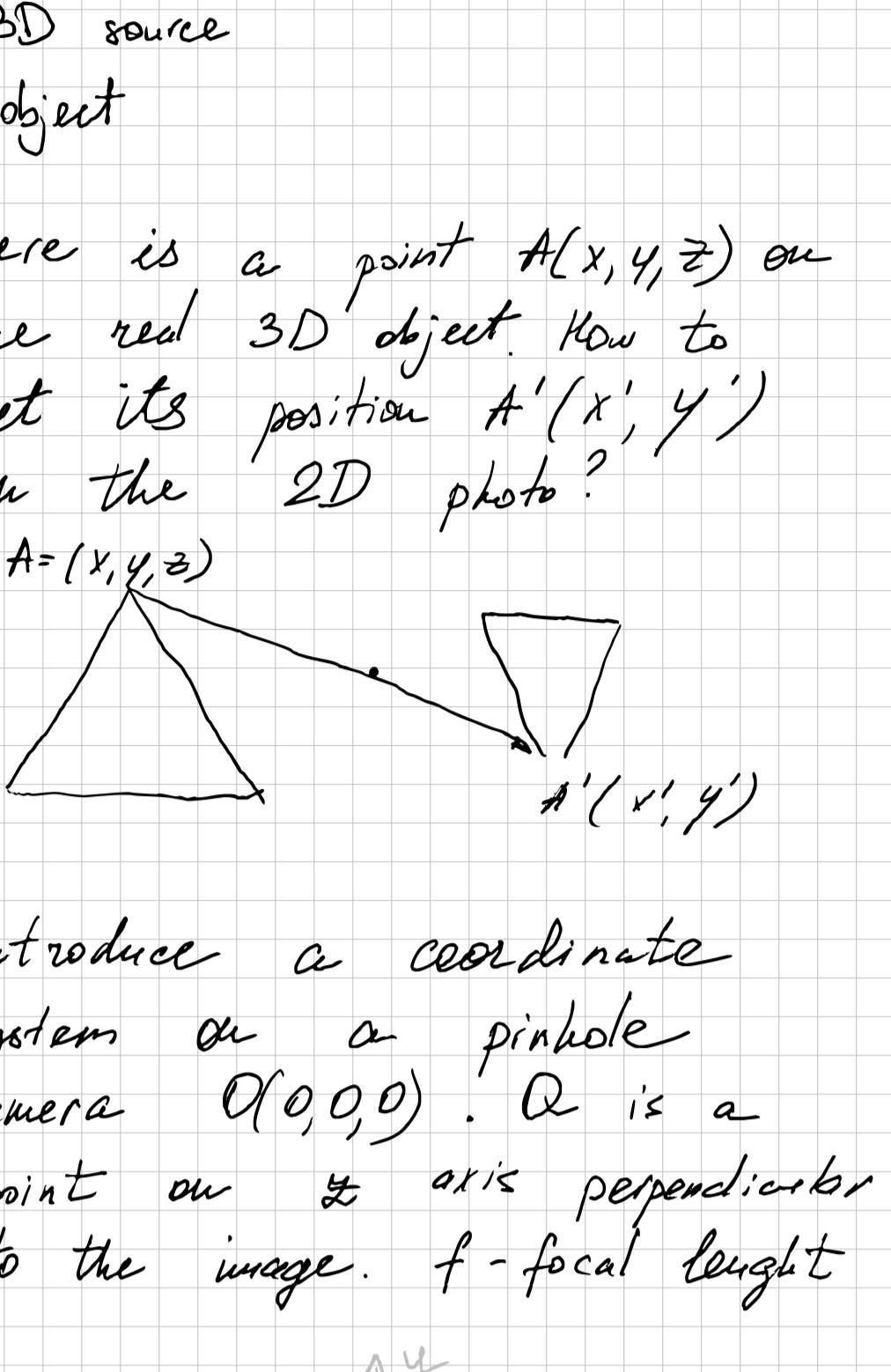
Solution: limit aperture width, so it does not let so many rays into sensor.

Assuming continuous space and image there will be more than one ray affecting infinitely many points of the image.

Solution: make aperture a single point with width 0.

Pinhole camera model

The diagram illustrates the pinhole camera model. On the left, a triangular shape represents a 3D scene. Two light rays originate from the top-left corner of this triangle and pass through a single point labeled "Aperture". From this aperture, the rays converge and form a smaller triangle on the right, labeled "2D image".



By symmetry around O:

$$x' = -C \cdot x; \quad y' = -C \cdot y$$

By similarity $\triangle APO \sim \triangle A'QO$: $\frac{QO}{PO} = \frac{f}{z}$

$\gamma(0,0,z)$ \rightarrow $A'(x', y', f)$

$$x' = \frac{f}{z} x \quad y' = \frac{f}{z} y$$

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \frac{f}{z} \begin{pmatrix} x \\ y \end{pmatrix}$$

In homogeneous coordinates:

$$\begin{pmatrix} x' \\ y' \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{f}{z} x \\ \frac{f}{z} y \\ 1 \end{pmatrix} \sim \begin{pmatrix} x \\ y \\ \frac{z}{f} \end{pmatrix}$$

Camera matrix

As a linear transformation:

$$\begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} \sim \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & \frac{1}{f} & 0 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$

Camera matrix:

$$C = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & \frac{1}{f} & 0 \end{pmatrix} \sim \begin{pmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$$

A canonical camera is a camera with focal length of 1 unit: $f=1$. It has a camera matrix:

$$C = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} = (I|0)$$

Pinhole Camera Model

Positive features:

- + Provides with a convenient 3D to 2D mapping.

Disadvantages:

- impossible to build in real life
- limits amount of light which leads to very dim images

Long Exposure Time

Possible solution: longer exposure time for more light coming through the pinhole.

Core problem: blurry images

Sources of such blur:

1. Moving objects (in different time they emit light from

different relative positions)
2 Not stabilised/shaking camera

Lenses

Converge all light rays into the same point.

The diagram illustrates the optical properties of a lens. On the left, a triangular shape representing a 3D object is shown. Several light rays originate from different points on the object's surface and travel towards a central vertical lens. The lens is depicted as a thick, elongated trapezoid. The light rays passing through the lens converge, as indicated by arrows pointing towards a single focal point on the right side of the lens. This focal point is where a second triangular shape, representing a 2D image, is formed. The image is inverted relative to the object.

3D object Lens 2D image

Upsides:

- + Solves the problem of pinhole camera

Downsides:

- Does not converge rays from every point of the object
- Different distances might lead to fuzzy parts of the image (unfocused objects)
- Various geometric distortions
- Might introduce vignetting

Camera physical matrix

Usually, each single pixel of

the camera matrix registers only one of the colors a raw image has only R, G, or B in every pixel. The layout of R, G, B filters is called a color filter mosaic.

The Bayer color filter (CFM) is one of the most popular layouts.

The Bayer CFM

In the Bayer CFM for each square 2×2 pixels, there are 2 pixels of

G , one R , one B . A camera matrix then consists of such repeated 2×2 squares with the same pattern.



Demosaicing

To recreate a full colored image demosaicing algorithms are used. Simplest methods either interpolate colors from a set of close neighbours or use a whole

2×2 square to describe
a single fully colored pixel
(reduces image width & height)

One 2×2 patch of the shape $(2, 2, 1)$

One colored pixel of the shape $(1, 1, 3)$

The diagram illustrates the process of averaging a 2x2 input patch to produce a 1x1 output pixel. On the left, a 2x2 grid of colored pixels is shown with labels: top-left is 'G R', top-right is 'R', bottom-left is 'B G', and bottom-right is 'G'. An arrow points from this input to the right, where a 1x1 grid is shown with a single label 'R' above it. Below the 1x1 grid is the formula $\frac{G+R+B+G}{4}$, representing the average of all four input pixels.

Often some additional steps like gamma correction or automatic white balance might be applied by the camera before outputting.

Depth Estimation

Depth estimation is the task of estimating distance relative to camera for all objects and surfaces.

Applications:

- 3D scene reconstruction
- self-driving cars
- 3D object reconstruction
- Augmented Reality

Datasets:

- KITTI
- Cityscapes
- SYNTHIA
- Middlebury
- NYUv2

Metrics:

- Mean absolute relative error (abs-rel)

$$\text{abs-rel} = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{d}_i - d_i\|}{d_i}$$

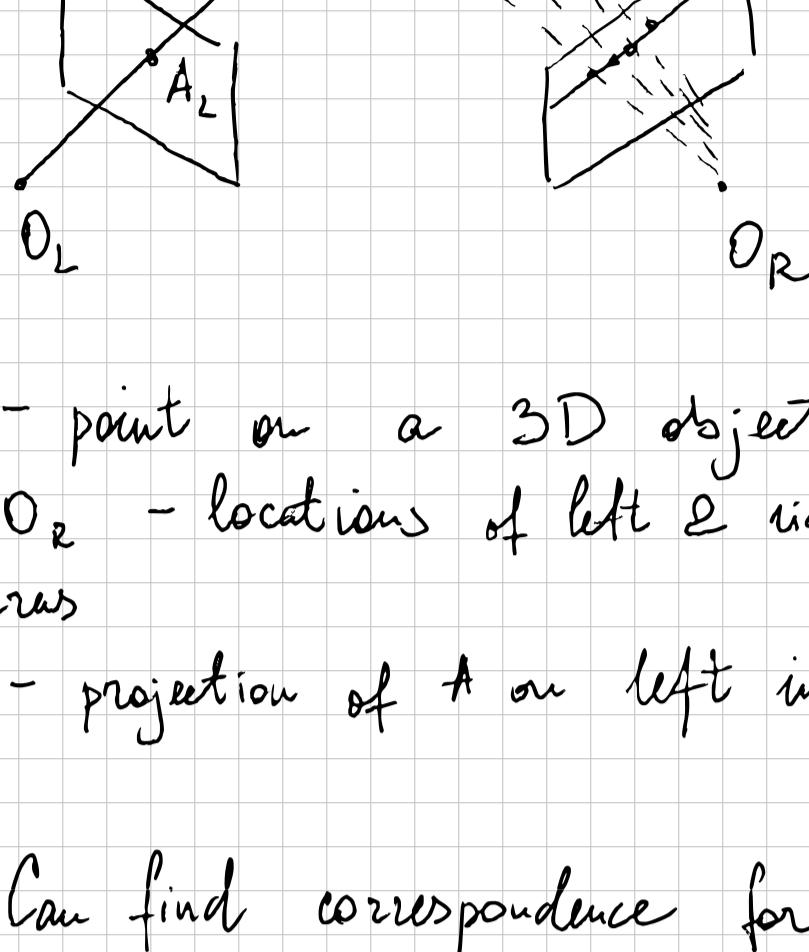
- Squared relative error

$$\text{sq-rel} = \frac{1}{N} \sum_{i=1}^N \frac{\|\hat{d}_i - d_i\|^2}{d_i}$$

- Root Mean Square Error

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N \|\hat{d}_i - d_i\|^2}$$

Stereo-based depth estimation

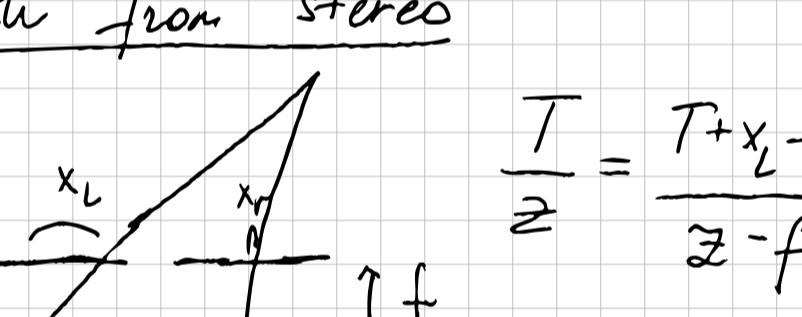


A - point on a 3D object

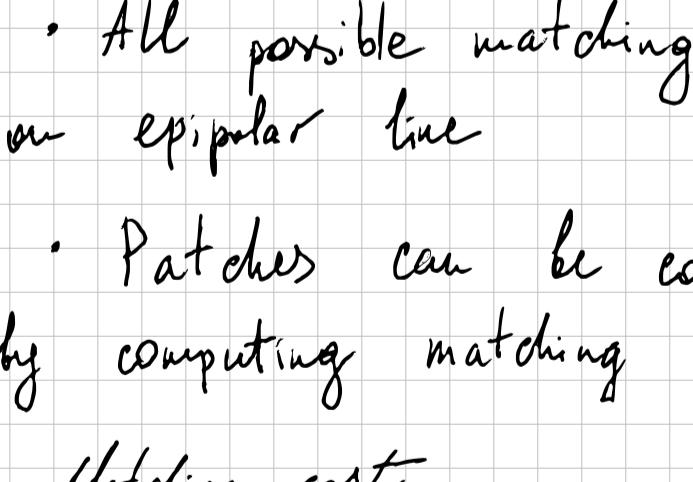
O_L, O_R - locations of left & right cameras

A_L - projection of A on left image

- Can find correspondence for A_L on O_R
- Can estimate 3D location of point A



Stereo rectification is a process of making parallel lines horizontal



How can we estimate disparity?

- All possible matching points on epipolar line
- Patches can be compared by computing matching cost

Matching cost

- Squared difference

$$(I_r - I_l)^2$$

- Normalised cross correlation

$$\frac{\sum [(I_l - \mu_{I_l})(I_r - \mu_{I_r})]}{\sqrt{I_l^2} \sqrt{I_r^2}}$$

- Deep matching cost

End-to-end disparity regression

- Learn deep representation

Monocular depth estimation

- Supervised
 - DORN
 - Introduce spacing - increasing discretization strategy
 - Remove subsampling in the last pooling layers
 - ASPP module consists of dilated convolutions
 - MIDAS
 - Use depth dataset obtained from 3D movies
 - Introduce scale and bias invariant loss
 - Train model on multiple dataset
 - Test on unseen data

Drawbacks

- Best models requires a lot of ground truth depth maps
- Hard to obtain high quality depth
- Time synchronization of a camera and depth sensor
- Poor generalization

Depth sources

- RADAR

Sparse, weather dependent

- LiDAR

Expensive, weather dependent

- 3D movies

Only relative depth

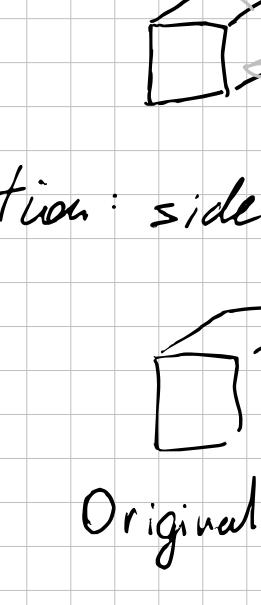
Self-supervised depth estimation

- Train require only stereo-pairs or videos
 - Model predicts depth and relative pose (6-DOF)

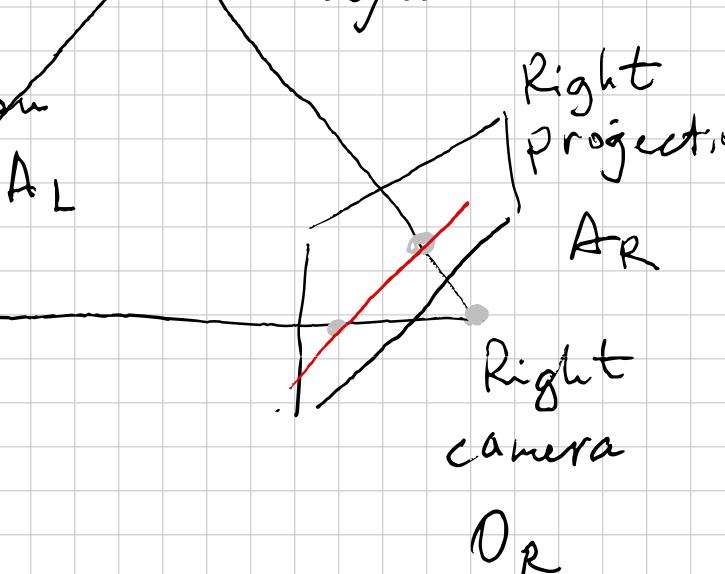
Monodepth 2

Epi polar geometry

cube or flat polyhedron?



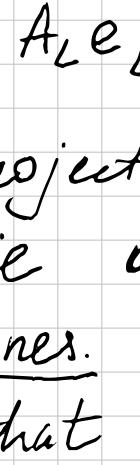
Issue: depth depends on flatness of an object. Information about it is lost due to 3D object \rightarrow 2D image translation.



Solution: side view

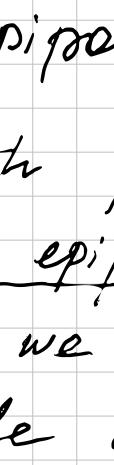


Original image



View from right

(3D)



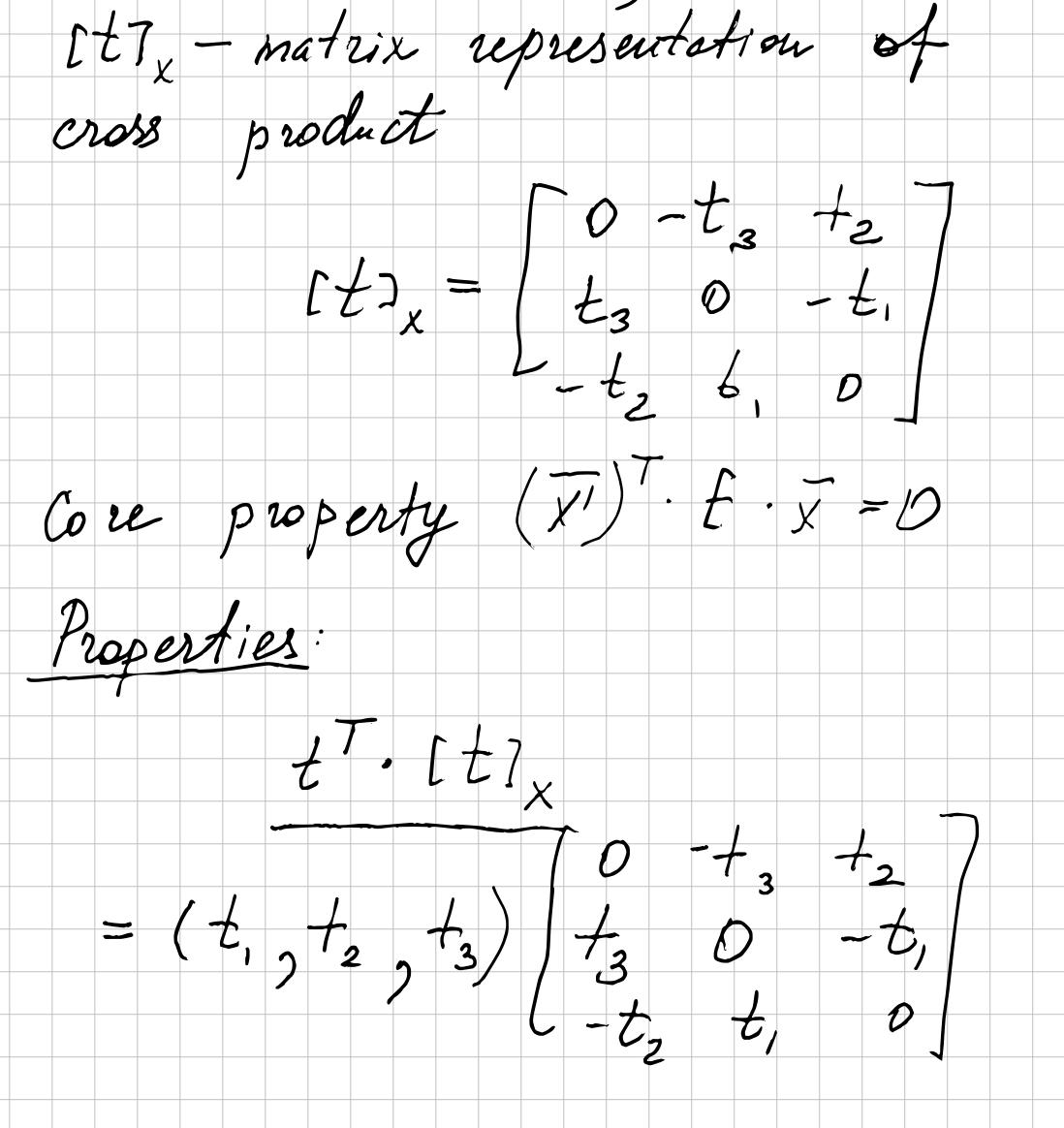
View from right

(flat object)

Epi polar geometry

2 cameras view the same 3D object from 2 different positions at the same time.

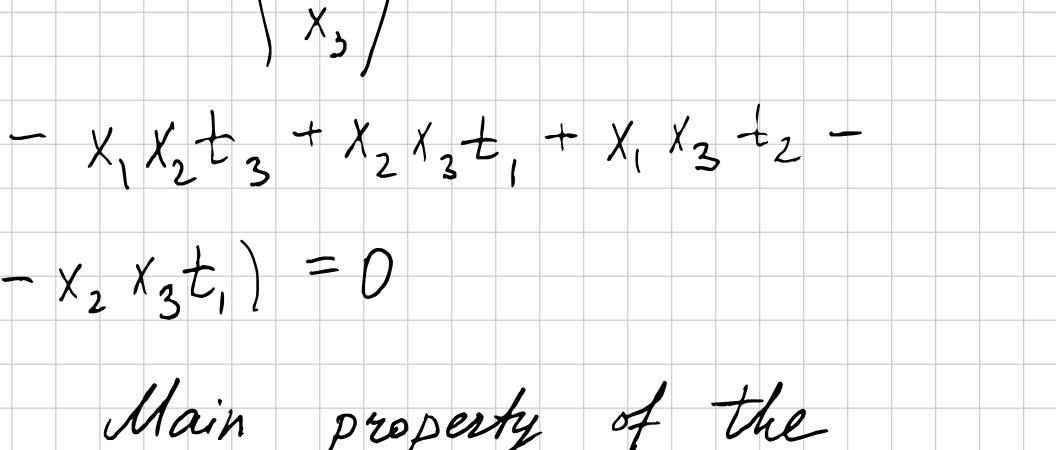
Result: 2 related 2D projections.



A - point on a 3D object

O_L, O_R - locations of left & right cameras

A_L, A_R - projections of A on left and right images



$O_L A D_R$ - epipolar plane

$O_L O_R$ - base line

e_L, e_R - left and right epipoles (intersection of a baseline and left/right image)

A_L, e_L, A_R, e_R - epipolar lines

Projections on both planes lie on respective epipolar lines. Knowing A_L we know that A_R will be on the right epipolar plane.

Consider coordinate systems related to 2 canonical cameras.

Point A is:

- left coordinate system: (a_1, a_2, a_3)
- right coordinate system: (a'_1, a'_2, a'_3)

Point A_L in left camera's coordinate system:

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = \frac{1}{a_3} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \Leftrightarrow \bar{x} = \frac{1}{a_3} \bar{a}$$

Point A_R in right camera's coordinate system:

$$\begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} = \frac{1}{a'_3} \begin{pmatrix} a'_1 \\ a'_2 \\ a'_3 \end{pmatrix} \Leftrightarrow \bar{x}' = \frac{1}{a'_3} \bar{a}'$$

A_R could be obtained from A_L by rotation and translation. Both are linear transformations.

$\bar{x}' = R \cdot (\bar{x} - t)$, where R - rotation; t - translation vector.

Essential matrix

$$\bar{x} = \frac{1}{a_3} \bar{a}$$

$$\bar{x}' = \frac{1}{a'_3} \bar{a}'$$

$$\bar{x}' = R(\bar{x} - t) \Leftrightarrow \bar{x}' = R \cdot \bar{x} - R \cdot t$$

$$R \cdot \bar{x} - R \cdot t = I \cdot \bar{x} - R \cdot t = \bar{x} - R \cdot t$$

$$(R \cdot \bar{x} - R \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

$$(R^T \cdot \bar{x} - R^T \cdot t)^T \cdot E \cdot \bar{x} = 0$$

Stereo systems

Stereo Pair of Images is a pair of images of the same scene made by different cameras in different positions.

- Stereo pair is used as an input or additional data (usually boosts model performance)
- Train with stereo pairs, inference with monocular distilled model
- Use sequences of images from a video as pseudo-stereo pair. Camera moves slightly between frames. With appropriate fps: not fast enough to change the scene, but fast enough to get a slightly different image

Localization and mapping

SLAM - simultaneous localization and mapping

localization - finding where the robot is with respect to the map

mapping - building a representation of the environment

Visual SLAM uses only images

Data sources

- Mono camera
- Stereo camera
- Depth sensor
- Inertial measurement unit (IMU)

Mono camera

- Low cost
- Scale ambiguity of a depth
- Depth can be estimated only by neural networks
- High computational cost

Stereo camera

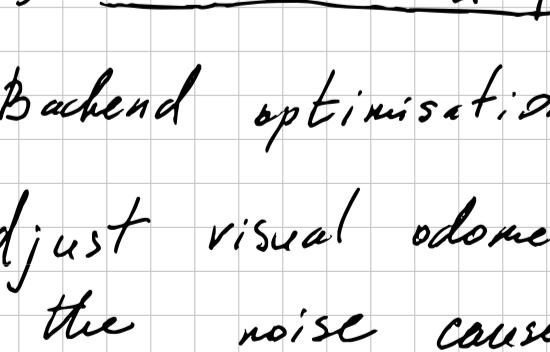
- Accurate absolute depth
- Complicated calibration process
- Very high computational cost

RGB-D camera

- True depth map
- Save computational resources
- Time synchronization with camera
- Noisy / sparse depth
- Issues with measure distance to transparent/reflecting surfaces

IMU

- Measures the angular velocity and the acceleration
- Accelerometer + gyroscope
- Measurements are noisy



Mathematical formulation

Input

$$U_{0:k} = \{u_0, u_1, \dots, u_k\}^T \text{ - input commands}$$

$$Z_{0:k} = \{z_0, z_1, \dots, z_k\}^T \text{ - sensors measurements}$$

Output

$$X_{0:k} = \{x_0, x_1, \dots, x_k\}^T \text{ - estimated positions}$$

$$m = \{m_0, m_1, \dots, m_k\}^T \text{ - landmarks}$$



1. Sensor data

Data is read from cameras

2. Visual odometry

Estimate the motion between consecutive camera frames.

Causes accumulative drift

3. Backend optimisation

Adjust visual odometry results for the noise caused by the sensors

4. Loop closing

Detect when a point has been identified twice, and correct the entire map for the accumulated drift error

5. Reconstruction

Build the map

Maps:

• Metric

• Topological

Problems of modern SLAM

• Lighting conditions

• Interacting with people

• Solution for joint indoor and outdoor environments

ORB SLAM

Parallel steps:

- Tracking
- Local mapping
- Loop closing
- Open source
- Real-time operations in large environments
- ORB feature detector
- Same features for all steps

ORB — oriented FAST and Rotated BRIEF

- FAST (Features from Accelerated Segment Test) is a high-speed corner detector
- BRIEF (Binary Robust Independent Elementary Features) is a feature point descriptor

FAST

- Checks 16 pixels in a neighborhood
 - Corner is detected if 12 sequential pixels are brighter or darker than the central pixel
 - First checks only 4 pixels : 1, 4, 5, and 13

16. 1. 2. 3.
· · ·
· P · · ·

BRIEF

- Smoothing an image
- Choosing random # of points for each image patch

• Comparison of the pixel intensity of pairs of points

• Result n-dimensional bit string

• n: 128, 256 or 512

Map points

- 3D position
- Viewing direction
- ORB descriptor
- Max 8 min distances at which the point can be observed

Key frames

- Camera pose
- Camera intrinsics : focal length, principal point
- ORB features

Tracking

1. ORB extraction
2. Initial Pose estimation from previous frame

3. If the tracking is lost: initial pose estimation via Global Relocalization

4. Track local map
5. New Key Frame decision

Local Mapping

Visibility graph :

- Graph of keyframes
- Edge if keyframes share map points

Essential graph :

- All keyframes
- Contains a spanning tree of edges

Loop Closing

1. Loop Candidates Detection
2. Compute the Similarity Transformation

3. Loop Fusion

4. Essential Graph Optimization

ORB SLAM-2

- Add stereo and 2D cameras

• More accurate

• Can estimate 3D position of a point

ORB SLAM-3

- Integrate IMU

• Improved recall-place recognition

Deep Learning for SLAM

- Visual odometry
- SLAM
- Feature detectors & descriptors

- DeepVO
- SSL VO + Depth
- CNN SLAM (monocular)
- Super Point

Novel view synthesis

Input: set of sparse images representing some object or scene, novel viewpoint position

Output: image representing object or scene from the given viewpoint

Approaches:

1. Modelling (object reconstruction)
2. Image-based rendering (NeRF)

Voxel-based approaches

3D-R2N2

Pix2Vox

TMV Net

Mesh-based approaches.

- Pixel2Mesh
- Textured Mesh Gen
- AtlasNet

Implicit shape representation

• Occupancy Networks