



# COMPUTER VISION

## EXERCISE 2 – STRUCTURE-FROM-MOTION AND STEREO

### 1 Pen and Paper

#### 1.1 Epipolar Geometry

- a) Assume you have two cameras, both with intrinsic parameters and rotations  $\mathbf{K} = \mathbf{R} = \mathbf{I}$ . Then for each of the three translation vectors  $\mathbf{t}_1$ ,  $\mathbf{t}_2$  and  $\mathbf{t}_3$  given below, compute the essential matrix, describe the orientation of the epipolar lines and determine location of the epiholes for the resulting camera configurations.

$$\mathbf{t}_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \mathbf{t}_2 = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}, \mathbf{t}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix},$$

- b) In the third case of the previous problem, where does the baseline lie and what does that imply for the location of the epiholes?
- c) When is the fundamental matrix equal to the essential matrix? Discuss your reasoning.

**Hint:** *Think about the relationship between camera- and image coordinates.*

#### 1.2 Triangulation

- a) Consider a system of two cameras with the following intrinsics and extrinsics:

$$\begin{aligned} \mathbf{K}_1 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, & \mathbf{K}_2 &= \begin{bmatrix} 2 & 0 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \\ \mathbf{R}_1 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, & \mathbf{R}_2 &= \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \\ \mathbf{t}_1 &= \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, & \mathbf{t}_2 &= \begin{pmatrix} -2 \\ 0 \\ 1 \end{pmatrix} \end{aligned}$$

Furthermore, assume that you have observations for a point in both cameras:

$$\tilde{\mathbf{x}}_1^s = \begin{pmatrix} 1/4 \\ 1/2 \\ 1 \end{pmatrix}, \tilde{\mathbf{x}}_2^s = \begin{pmatrix} -1/5 \\ 1/5 \\ 1 \end{pmatrix}$$

For the given system, triangulate the 3D point  $\tilde{\mathbf{x}}_w$  in world coordinates that corresponds to the observations. You can assume that the observations are exact.

#### 1.3 Stereo Vision

- a) Show that for a stereo camera system the depth measurement error grows quadratically with depth.
- b) You can also think of this relationship as the depth resolution of the stereo camera system. How can we change the system setup to get a better depth resolution? What disadvantages might this have?

## 1.4 Block Matching

- a) Consider two  $K \times K$  windows of pixels flattened to vectors  $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^{K^2}$ . Show that the *Zero Normalized Cross-Correlation (ZNCC)* is invariant to changes in brightness in these windows. For this you can assume changes in brightness to be linear transformations of the form  $\mathbf{w}'_i = \alpha_i \mathbf{w}_i + \mathbf{1}\beta_i$ , where  $\mathbf{1} \in \mathbb{R}^{K^2}$  is a vector of all ones.
- b) You are given the following pair of  $5 \times 7$  stereo images, where part of the background is occluded by a thin structure (represented by the column of 10s). Using the *Sum of Squared Differences (SSD)* with a  $3 \times 3$  window, determine the *Winner-Takes-All (WTA)* disparity  $d \in \{0, 1, 2\}$  for the background pixel marked in Cyan. Discuss your result.

|   |   |   |    |   |   |   |
|---|---|---|----|---|---|---|
| 1 | 2 | 3 | 10 | 5 | 6 | 7 |
| 1 | 2 | 3 | 10 | 5 | 6 | 7 |
| 1 | 2 | 3 | 10 | 5 | 6 | 7 |
| 1 | 2 | 3 | 10 | 5 | 6 | 7 |
| 1 | 2 | 3 | 10 | 5 | 6 | 7 |

(a) Left image

|   |    |   |   |   |   |   |
|---|----|---|---|---|---|---|
| 2 | 10 | 4 | 5 | 6 | 7 | 8 |
| 2 | 10 | 4 | 5 | 6 | 7 | 8 |
| 2 | 10 | 4 | 5 | 6 | 7 | 8 |
| 2 | 10 | 4 | 5 | 6 | 7 | 8 |
| 2 | 10 | 4 | 5 | 6 | 7 | 8 |

(b) Right image

- c) Below we have the full disparity maps for the images from the previous exercise computed from the left- to the right image and from the right- to the left image, respectively. Perform a left-right consistency test for the pixels marked in Cyan, Green and Red. Which of the points pass the test? Is the test succesful in determining incorrect disparity estimates?

**Remark:** The disparities were computed in the same way as in the previous exercise. To compute disparities along the image boundaries, the images were padded with zeros.

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 2 | 2 | 1 | 0 |
| 0 | 1 | 2 | 2 | 2 | 1 | 0 |
| 0 | 1 | 2 | 2 | 2 | 1 | 0 |
| 0 | 1 | 2 | 2 | 2 | 1 | 0 |
| 0 | 1 | 2 | 2 | 2 | 1 | 0 |

(a) Left  $\rightarrow$  Right

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| 2 | 2 | 2 | 2 | 1 | 0 | 0 |
| 2 | 2 | 2 | 2 | 1 | 0 | 0 |
| 2 | 2 | 2 | 2 | 1 | 0 | 0 |
| 2 | 2 | 2 | 2 | 1 | 0 | 0 |
| 2 | 2 | 2 | 2 | 1 | 0 | 0 |

(b) Right  $\rightarrow$  Left

## 1.5 Learned Stereo and End-to-End Models

- a) Recent approaches in end-to-end disparity estimation often build a disparity cost volume and apply 3D convolutions to it to estimate the final disparity map. However, 3D convolutions are computationally expensive, limiting both the resolution and maximum disparity that can be used.

Below we consider two sequences of two 2D- and 3D convolutions, respectively, applied to the same input tensor. We describe the layer configurations as  $\text{ConvND}(C_{in}, C_{out}, k)$ , where input channels are denoted by  $C_{in}$ , output channels by  $C_{out}$  and  $k$  is the kernel size. For simplicity you can assume square kernels, as well as appropriate padding and a stride of one, such that the spatial dimensions remain unchanged. Shapes are specified by the number of channels followed by the spatial dimensions, so (Channels, Height, Width) for 2D convolutions and (Channels, Depth, Height, Width) for 3D convolutions.

For both sequences, calculate the total amount of memory required to store the activations and trainable parameters. For this you can fill in the blank fields in the table below.

| Layer             | Input Shape       | Output Shape | # Trainable Parameters | Memory |
|-------------------|-------------------|--------------|------------------------|--------|
| Conv2D(32, 64, 3) | (32, 128, 128)    |              |                        |        |
| Conv2D(?, 128, 3) |                   |              |                        |        |
| Conv3D(1, 64, 3)  | (1, 32, 128, 128) |              |                        |        |
| Conv3D(?, 128, 3) |                   |              |                        |        |

- b) You are working with a GC-Net-style architecture to solve a disparity estimation problem. For two particular pixels  $p_1$  and  $p_2$  in the cost volume, the network estimates the following matching costs:

- $p_1$ :  $c_\theta(d) = [1.0, 3.0, 10.0, 3.0, 1.0]$
- $p_2$ :  $c_\theta(d) = [10.0, 2.0, 1.0, 2.0, 10.0]$

where  $d \in \{0, 1, 2, 3, 4\}$ . For both pixels, calculate the expectation of the disparity and discuss the result. GC-Net uses the discrepancy between the expected and the ground-truth disparity as a loss function during training. What kind of behaviour does this encourage?

**Hint:** Compare the distributions over disparities implied by the cost vectors with the final result.

## 2 Coding Exercises

The following coding exercises are split into two sections: A structure-from-motion and a stereo section, with corresponding jupyter notebooks `code/sfm/sfm.ipynb` and `code/stereo/stereo.ipynb`. As in the last exercise, the notebooks are self-contained but you can also use this document as guidance. If you are stuck, you can find *Hints* in the notebooks themselves which are written upside-down.

### 2.1 Structure-From-Motion

- Implement the function `compute_fundamental_matrix` which takes sets of corresponding keypoints  $\bar{\mathbf{x}}_i$  from the first and second image as input and returns the fundamental matrix  $\mathbf{F}$  using the 8-point algorithm.
- Implement the function `compute_fundamental_matrix_normalized` which again takes sets of corresponding keypoints  $\bar{\mathbf{x}}_i$  from the first and second image as input and returns the fundamental matrix  $\tilde{\mathbf{F}}$ , but this time using the *normalized* 8-point algorithm.
- Implement the function `compute_essential_matrix` which takes the fundamental matrix  $\tilde{\mathbf{F}}$  as well as the intrinsics  $\mathbf{K}_i$  for the first and second image as input and returns the essential matrix  $\tilde{\mathbf{E}}$ .
- Implement the function `triangulate_point`. This function takes keypoints  $\bar{\mathbf{x}}_i$ , intrinsics  $\mathbf{K}_i$  as well as the relative rotation  $\mathbf{R}$  and translation  $\mathbf{t}$  from two images as input and returns a triangulated 3D point  $\tilde{\mathbf{x}}_w$ .

### 2.2 Stereo

- Implement the function `sad`, which given a window size and maximum disparity  $D$ , takes a stereo image pair as input and returns a disparity map, computed from the left to the right image. If you are interested you can also implement the bonus function `sad_convolve`, which should be significantly faster than `sad`.
- Create a visualization of the computed disparities by implementing the function `visualize_disparity`. It's a good idea to also visualize the input images to see if the results are sensible.
- Experiment with different window sizes (for example 3, 7, 15) and report which one leads to better visual results and why? In case you were not able to solve the previous exercises, you can use the provided disparity maps in the `code/stereo/examples/` folder.
- Why do you think the block matching approach fails to lead to good estimations around homogeneous regions such as the road?
- Develop a Siamese Neural Network architecture. For this you can use the `StereoMatchingNetwork` class. In particular, implement the `__init__` and `forward` methods to initialize the layers and define the forward pass. Details on the architecture can be found in the jupyter notebook.

- f) From the lecture you know two classes of Siamese Neural Network architectures - which class does the architecture you implemented in **StereoMatchingNetwork** belong to?
- g) Implement the function `calculate_similarity_score`, which takes an instance of the Siamese Neural Network as well as patches  $\mathbf{w}_L$  and  $\mathbf{w}_R$  from the left and right image as input and returns the similarity of those patches. The similarity should be computed based on features extracted by the Siamese Network.
- h) Try to improve the network by finding better hyperparameters. For example, you can vary the number of training iterations or the number of filters in the convolutional layers. Explain your findings.
- i) Compare the visualization of the disparity maps from the Siamese Neural Network to the ones obtained by the block matching algorithm. Which predictions are better and why? Can you find regions in the scenes where the differences in predictions are most dominant? (If you were not able to solve the previous exercises, you can use the provided disparity maps in the `code/stereo/examples/` folder.)