

Computer Vision Notes

Huu Duc Nguyen M.Sc.

29 March 2022

Contents

| | |
|---|-----------|
| Abbreviations | 1 |
| 1 Introduction | 2 |
| 2 Mathematics Backgrounds | 3 |
| 2.1 The Matrix Equation | 3 |
| 3 Image Formation | 4 |
| 3.1 Camera Obscura | 4 |
| 3.2 Pinhole Camera | 4 |
| 3.3 Digital image | 5 |
| 3.4 Color Sensing | 5 |
| 3.4.1 Demosaicing | 5 |
| 4 Image Processing | 7 |
| 4.1 Linear Filters | 7 |
| 4.2 Background | 8 |
| 4.3 Non-Linear Filters | 9 |
| 4.4 Multi-Scale Representations | 9 |
| 4.5 Filters as Templates | 9 |
| 4.6 Image Gradients | 10 |
| 4.7 Edge Detection | 10 |
| 4.8 Structure Extraction | 10 |
| 5 Segmentation | 11 |
| 6 Object Detection | 12 |
| 7 Local Feature | 13 |
| 8 Deep Learning for CV | 14 |
| 9 3D Computer Vision | 15 |
| 9.1 Introduction | 15 |
| 9.2 Depth Extraction | 15 |
| 9.3 3D Shape representation | 16 |
| 9.3.1 Voxel Grid | 16 |

| | | |
|-----------|---|-----------|
| 9.3.2 | Point Cloud | 16 |
| 9.3.3 | Mesh | 16 |
| 9.3.4 | Occupancy | 16 |
| 9.4 | Classic 3D Reconstruction | 16 |
| 9.4.1 | Epipolar Geometry | 17 |
| 9.4.2 | Stereo Image Rectification | 17 |
| 9.4.3 | Correspondence Search | 17 |
| 9.4.4 | Stereo Reconstruction | 18 |
| 9.4.5 | Camera Calibration | 18 |
| 9.4.6 | Eight Point Algorithm | 18 |
| 9.5 | Deep Learning for 3D CV | 18 |
| 10 | Single Object Tracking | 19 |
| 11 | Bayesian Filtering | 20 |
| 12 | Multi Object Tracking | 21 |
| 13 | Visual Odometry | 22 |
| 14 | SLAM | 23 |
| 15 | Deep Learning for Video Analysis | 24 |
| | Bibliography | I |

Abbreviations

| | |
|---------------|---------------------------------------|
| info. | information |
| a.k.a. | also known as |
| no. | number of |
| func. | function |
| vs. | versus |
| i.i.d. | independent & identically distributed |
| LSI | linear shift invariant |
| SVD | Singular Value Decomposition |

1 Introduction

Goal: The goal of computer vision is enabling machine to understand images & videos.
There are two major tasks:

- measurement: compute properties of 3D world (distance, shape)
- perception & interpretation: recognize objects, people, activities, ..

Outlines

- Chap. 2 presents mathematics backgrounds for computer vision
- [TODO: Chap. ?? do sth]

2 Mathematics Backgrounds

This chapter presents some mathematics backgrounds for later problems.

2.1 The Matrix Equation

Problem: Solve $Ax = 0$

- Applying Singular Value Decomposition ([SVD](#)) for matrix A

$$A = U.D.V^T = U. \begin{bmatrix} d_{11} & \cdots & d_{1N} \\ \vdots & \ddots & \vdots \\ d_{N1} & \cdots & d_{NN} \end{bmatrix} \cdot \begin{bmatrix} v_{11} & \cdots & v_{1N} \\ \vdots & \ddots & \vdots \\ v_{N1} & \cdots & v_{NN} \end{bmatrix}^T$$

- Solution of $Ax = 0$ is the null space vector of A , which corresponds to the smallest (last) singular vector of A : $[v_{1N}, \cdots, v_{NN}]^T$.

3 Image Formation

3.1 Camera Obscura

also known as (a.k.a.) the "Dark Chamber" (Leonardo Da Vinci, 1545)

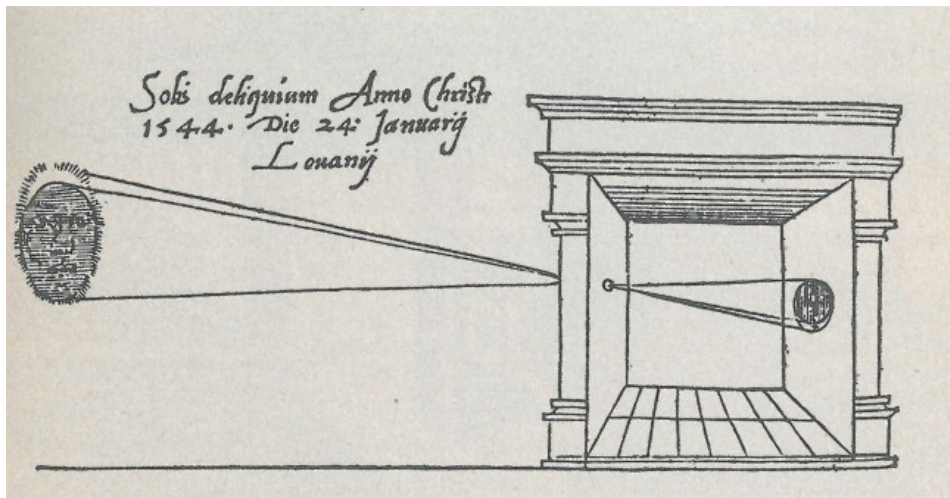


Figure 3.1: Camera obscura [Fri45].

3.2 Pinhole Camera

- Pinhole size = aperture
 - too big \Rightarrow blurring
 - too small \Rightarrow also blur, but because of diffraction
 - but then, *image is dark*
- \Rightarrow Use lenses: keep image sharp while *capture more light*
- The thin lens
- Focus & Depth of Field:
 - Large aperture: small depth of field
(only object within the correct distance will be at focus, while background is blur)
 - Small aperture: large depth of field, but need more light
- The lens focus $f \gtrless$ field of view
 - f gets smaller \Rightarrow wide-range image
 - f gets greater \Rightarrow telescopic image

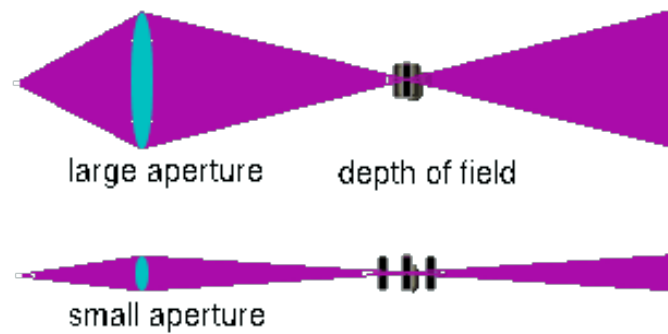


Figure 3.2: Varied depths of field depending on aperture size.

3.3 Digital image

- Discretize the image into a grid of pixels
- Quantize light intensities \Rightarrow pixel values
- Resolution: number of (no.) pixels (most commonly understand)

3.4 Color Sensing

Referring to the process of assigning pixel values from color information of world objects.

- Color image: RGB is just 1 of many color spaces, e.g., LUV, XYZ ([Wikipedia](#)).
- Grey-scale image

3.4.1 Demosaicing

Digital camera takes in light through a filter (Bayer or Xtrans) \Rightarrow we get a gray-scale image (Fig. 3.3). We need to apply demosaicing based on the filter's pattern to get the color image from the raw image. Sources: [YouTube](#), [Wikipedia](#).

NOTE: Raw image has a ***green cast***

Twice many green as red & blue, because human eyes are twice as sensitive to the green part to other red or blue part.

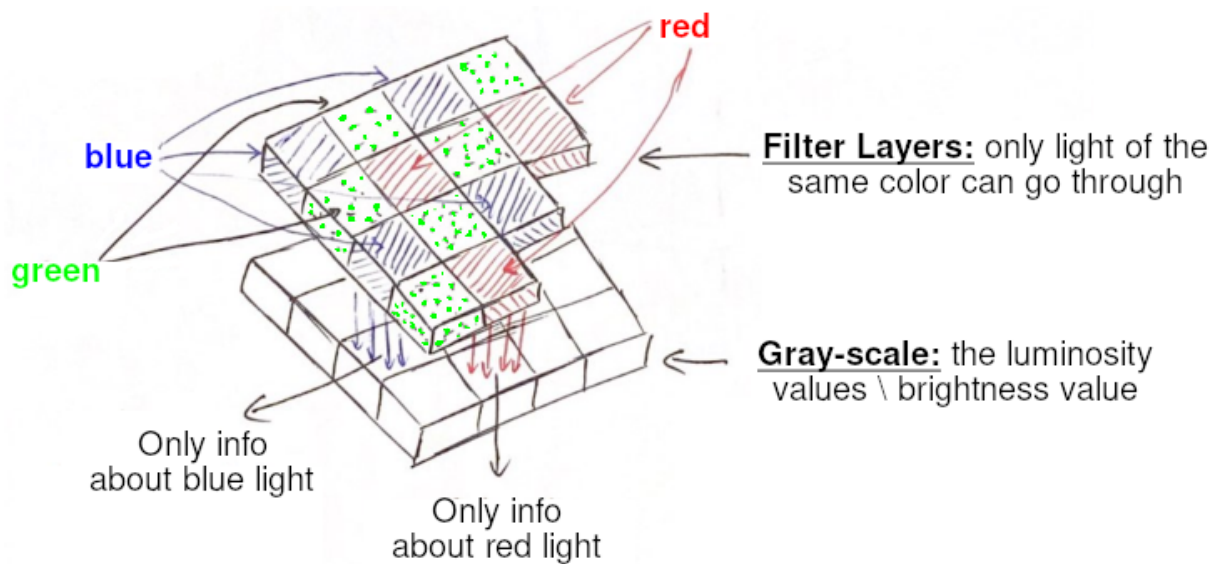


Figure 3.3: E.g. Bayer Filter. In the raw image , which lies below the filter layers, each pixel only has information (info.) of only 1 among 3 light sources. Demosaicing uses the values of surrounding pixels to infer the brightness of other light sources.

4 Image Processing

4.1 Linear Filters

Types of noise:

- Salt & pepper noise
- Impulse noise
- Gaussian noise

$$noise = randn(size(img)) \times \sigma$$

$$output = img + noise$$

- **Basic assumption:** independent & identically distributed (i.i.d.)

Types of filter:

- Correlation Filter: $G[i, j] = \frac{1}{(2k+1)^2} \sum_{u=-k}^k \sum_{v=-k}^k F[i+u, j+v]$

$$\text{different weights: } G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i+u, j+v] \Rightarrow \boxed{G = H \otimes F}$$

with $H[u, v]$ as non-uniform weights

Matlab: filter2, imfilter

- Convolution: $G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i-u, j-v] \Rightarrow \boxed{G = H * F}$

Matlab: conv2

If $H[u, v] = H[-u, -v] \Rightarrow$ correlation \equiv convolution

- Averaging Filter: **Ringing Artifacts??**

- Gaussian Filter: $\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

Rule of thumb: set the filter width to 6σ

More noise $\Rightarrow \uparrow \sigma \Rightarrow$ blurring effect

NOTE:

- **k is from the window size $(2k+1) \times (2k+1)$**
- **Efficient implementation:** if filter is separable \Rightarrow apply 1D filter 2 times to have a 2D filter \Rightarrow Reduce the computational cost from $\mathcal{O}(K^2)$ to $\mathcal{O}(2K)$, with K as the kernel size
- When coding with Python, the origin of image plane is top left corner, x -axis goes left, y -axis goes downward (Fig. 4.1)

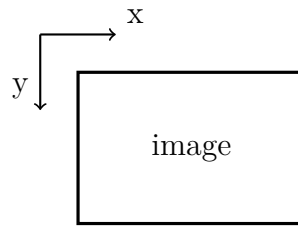


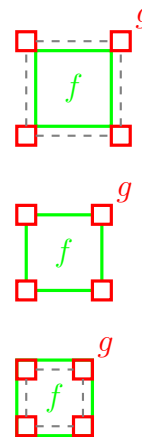
Figure 4.1: Image coordinate system in Python

- Boundary issues:

- Full: output size = $f + g$

- Same: output size = f

- Valid: output size = $f - g$



Pixel near boundary:

- Clip filter (black) \Rightarrow dark border
- Wrap around
- Copy edge \Rightarrow Strong edge response
- Reflect across edge
- Correlation versus (vs.) convolution:
 - Both are linear shift invariant linear shift invariant (LSI):

$$h \circ (f_0 + f_1) = h \circ f_1 + h \circ f_0$$
 - Conv is better, it has additional nice properties
 - * commutative: $f * g = g * f$
 - * associative: $(f * g) * h = f * (g * h)$
 - * Fourier transform $f * g \rightarrow F.G$ and $f.h \rightarrow F * H$
 - With impulse image, Conv reproduces itself, while Corr reflects itself.

4.2 Background

- Taking the Fourier Transform of a signal \Rightarrow Frequency coefficients \Rightarrow **Frequency Spectrum**
- Duality:** The **better** a function is **localized** in one domain the **worse** it is **localized** in the other domain.

- Effect of Convolution: $f * g \rightarrow F \cdot G$
 taking convolution in one domain is equivalent to multiplication in the other domain
 A Gaussian has compact support in both domains
 \Rightarrow convenient choice for **low-pass filter**
- Sharpening filter - **high-pass filter**: emphasizes noise as well

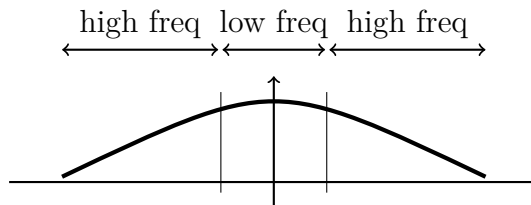


Figure 4.2: Frequency domain (Fourier).

4.3 Non-Linear Filters

- Median filter: replace each pixel by the median of the neighbors.
 - **remove spikes** (good for impulse, salt & pepper noise)
 - **edge preserving** (unlike mean filter)

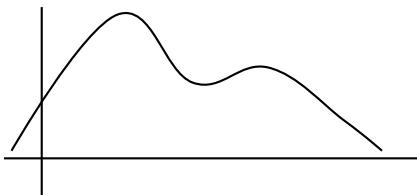
NOTE: If we increase the Median filter's filter size \Rightarrow reduce structure and loose details

4.4 Multi-Scale Representations

Image pyramid: very little overhead (in terms of computational cost).

Fourier Interpretation: Discrete Sampling

Sampling in spatial domain is like multiplying with a spike function (func.).



\Rightarrow Sampling in the frequency domain is like convolving with a spike func.

4.5 Filters as Templates

Correlation filtering as Template Matching.

4.6 Image Gradients

4.7 Edge Detection

4.8 Structure Extraction

[TODO: missing content]

5 Segmentation

[TODO:]

6 Object Detection

[TODO:]

7 Local Feature

[TODO:]

8 Deep Learning for CV

[TODO:]

9 3D Computer Vision

9.1 Introduction

3D Computer Vision gives a representation that is closer of things that we interact in our lives. Thus, it will empower various novel applications in:

- Autonomous Driving
- Robotics
- Remote Sensing
- Medical Treatment
- Design Industry
- Augmented Reality

[**TODO:**] Learning resources: [??](#).

3D computer vision problems includes:

- Depth extraction
- 3D Reconstruction
- Object Classification
- Object Detection
- Object Segmentation
- ??

Challenges of 3D computer vision:

- something here

9.2 Depth Extraction

The goal: extract the depth, as the 3rd dimension for a 2D image.

The depth map is a simple grey image with values in range $[0, 255]$, 0 for point afar and 255 for points in near distances.

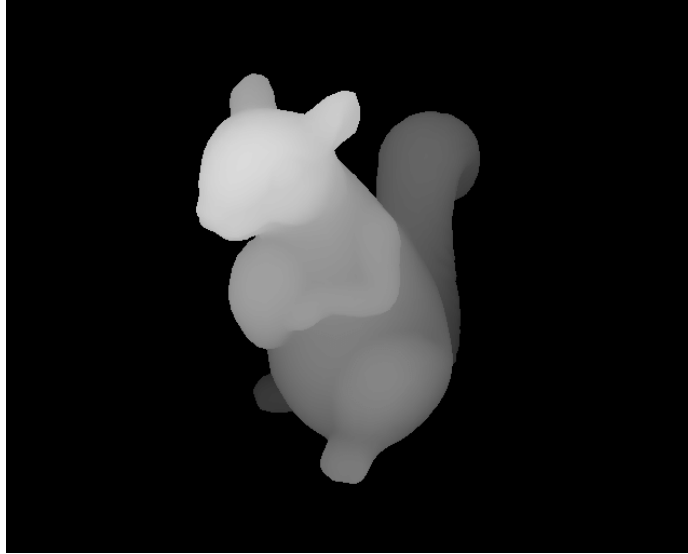


Figure 9.1: Example of a depth map [TSS+18].

9.3 3D Shape representation

There are explicit representations and implicit representations, where parametric functions are used to differentiate a specific point is inside or outside the shape, or the distance to the shape surface. Typically, the parametric functions are in form of neural networks

9.3.1 Voxel Grid

9.3.2 Point Cloud

9.3.3 Mesh

9.3.4 Occupancy

9.4 Classic 3D Reconstruction

Geometric vision:

- Visual Cues (Details)
 - Shading
 - Texture
 - Focus
 - Perspective
 - Motion

- Stereo vision: process of extracting 3D information from multiple 2D views of a scene

9.4.1 Epipolar Geometry

Epipolar geometry is the geometry of stereo vision. The **basic principle** of epipolar geometry is **triangulation** of points. In Fig. 9.2, O_1 and O_2 are the camera poses, X_1

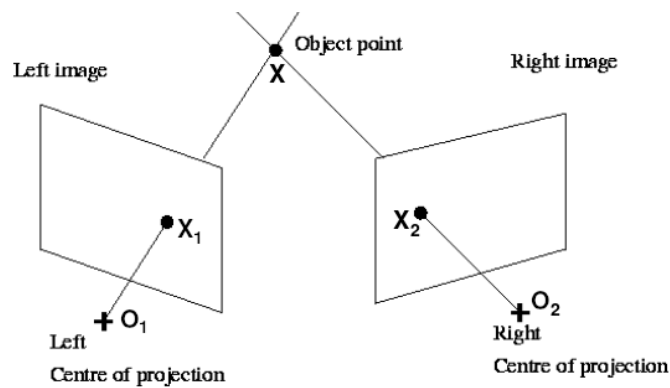


Figure 9.2: Example of triangulation ([src](#)). The lines connecting the camera poses with the correspondent points must intersect at the real object world space.

and X_2 are the correspondent points on each image planes, and X is the real object point in world space.

[TODO:]

9.4.2 Stereo Image Rectification

Re-project image planes on to a common plane, which is parallel to the baseline
 \Rightarrow Scan lines are epipolar lines.

[TODO: Add images]

9.4.3 Correspondence Search

Correspondence search simple means matching a point with another point in a different image.

NOTE: In practice, use both.

| Dense Correspondence Search | Sparse Correspondence Search |
|--|--|
| <ul style="list-style-type: none"> • For each pixel, find correspondence • Easy when epipolar lines are scan lines (apply rectification) | <ul style="list-style-type: none"> • Only for a set of detected feature • Use feature description (Harris, SIFT??) |
| — Pros — | |
| <ul style="list-style-type: none"> • Simple process • More depth \Rightarrow useful for surface re-construction | <ul style="list-style-type: none"> • Efficiency • Can have more reliable matches • Less sensitive to illumination \Rightarrow robust |
| — Cons — | |
| <p>Problem with:</p> <ul style="list-style-type: none"> • texture-less regions • different viewpoints | <ul style="list-style-type: none"> • Have to know enough to pick good features • Sparse information |

9.4.4 Stereo Reconstruction

Main steps:

- Calibrate cameras
- Rectify images
- Compute disparity
- Estimate depth

This is just the ideal case.

- What if, how can we get extrinsic **info.** from calibration?
- What to do when triangulation failed?

9.4.5 Camera Calibration

9.4.6 Eight Point Algorithm

9.5 Deep Learning for 3D CV

10 Single Object Tracking

[TODO:]

11 Bayesian Filtering

[TODO:]

12 Multi Object Tracking

[TODO:]

13 Visual Odometry

[TODO:]

14 SLAM

[TODO:]

15 Deep Learning for Video Analysis

[TODO:]

Bibliography

- [Fri45] R. G. Frisius. *De radio astronomico et geometrico liber*. Ap. Gul Cavellat, 1545.
- [TSS+18] H. Tjaden, U. Schwanecke, E. Schömer, and D. Cremers. “A Gauss-Newton Approach to Real-Time Monocular Multiple Object Tracking”. In: (July 2018).