

Pinhole camera

and homogeneous coordinates

Morten R. Hannemose, mohan@dtu.dk

February 2, 2024

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



Learning objectives

After this lecture you should be able to:

- explain homogeneous coordinates
- convert to and from homogeneous coordinates
- perform relevant coordinate transformations
- explain the pinhole camera model

Presentation topics

Pinhole camera

Perspective transformations

Rigid transformations

Homogeneous coordinates

Lines in homogenous coordinates

Summary of homogeneous coordinates

Pinhole camera model

Intrinsics

Extrinsics

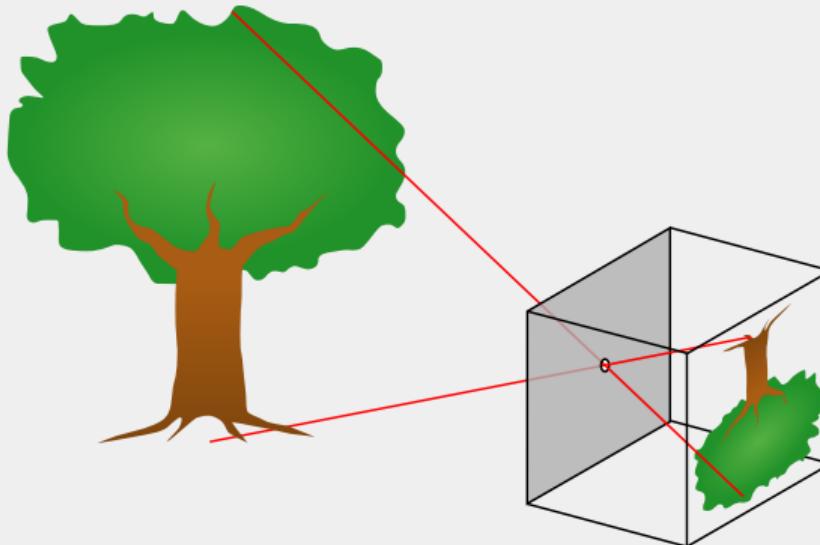
Pinhole camera

What is a “good” camera model?

... in terms of accuracy vs. ease-of-use?

1. As accurate as possible
2. As easy to use as possible
3. somewhere between 1 and 2

Pinhole camera



Light travels in straight lines.

The projected image appears upside down.

**Each point in an image corresponds to a
direction from the camera**

Real life example

Can I get two volunteers?

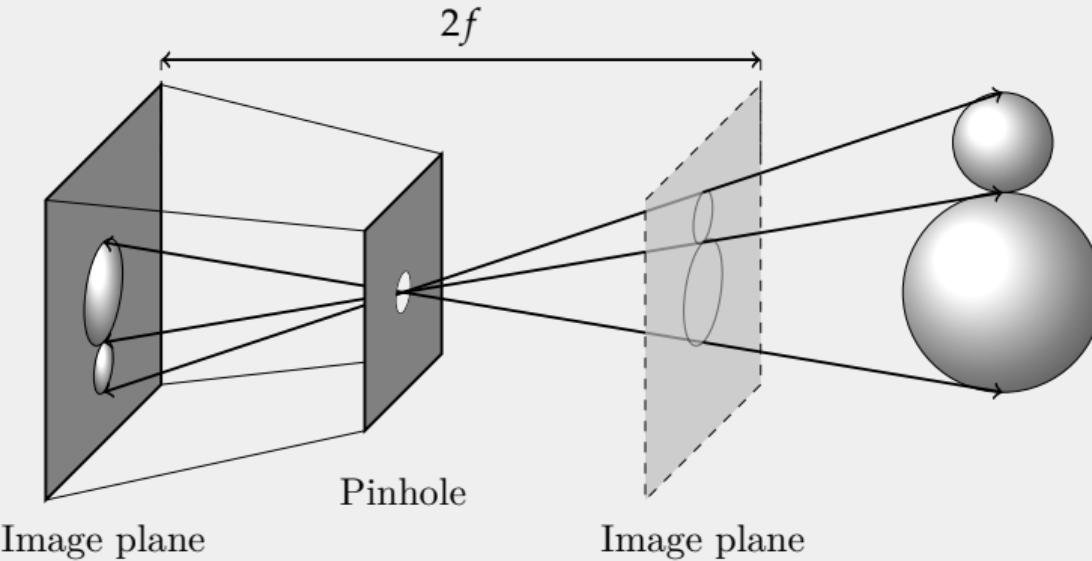
Real life example

Can I get two volunteers?

- A point seen in a single camera must be along a specific line in the other camera.
- Seeing the same point in two cameras is enough to find the point in 3D.

Perspective transformations

Pinhole camera again

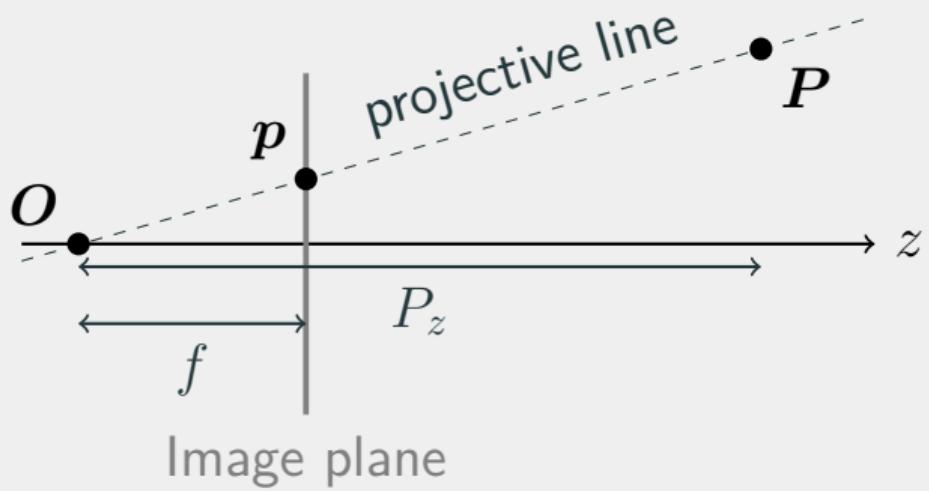


When modelling we place the “image plane” in front of the camera.
The distance from image plane to camera is the **focal length** f .

Perspective projection

$$\mathbf{P} = \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix},$$

$$\mathbf{p} = \frac{f}{P_z} \begin{bmatrix} P_x \\ P_y \end{bmatrix}$$

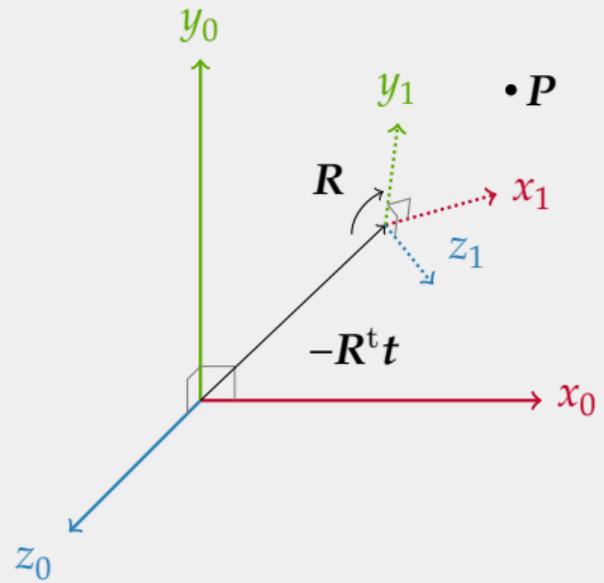


Rigid transformations

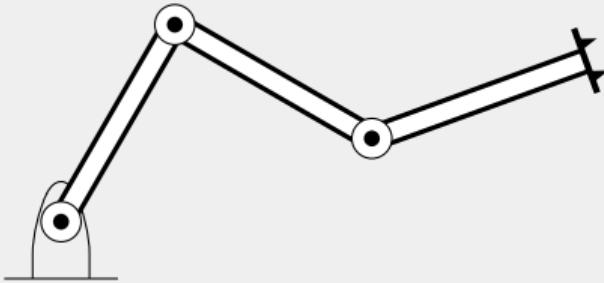
Rigid transformations: rotations and translations

Rotation matrix R and translation vector t

$$P_1 = RP_0 + t$$



Robot arm transformations



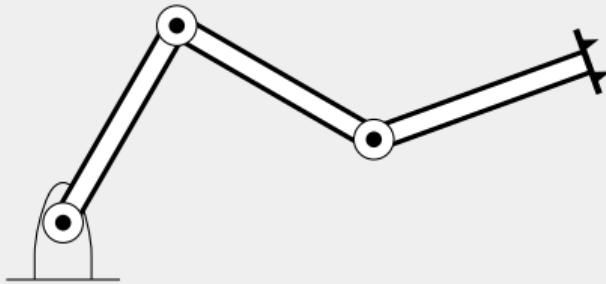
$$\mathbf{P}_1 = \mathbf{R}_1 \mathbf{P}_0 + \mathbf{t}_1$$

$$\mathbf{P}_2 = \mathbf{R}_2 \mathbf{P}_1 + \mathbf{t}_2$$

$$\mathbf{P}_3 = \mathbf{R}_3 \mathbf{P}_2 + \mathbf{t}_3$$

$$\mathbf{P}_4 = \mathbf{R}_4 \mathbf{P}_3 + \mathbf{t}_4$$

Robot arm transformations



$$\mathbf{P}_4 = \mathbf{R}_4(\mathbf{R}_3(\mathbf{R}_2(\mathbf{R}_1\mathbf{P}_0 + \mathbf{t}_1) + \mathbf{t}_2) + \mathbf{t}_3) + \mathbf{t}_4$$

$$\mathbf{P}_4 = \mathbf{R}_4\mathbf{R}_3\mathbf{R}_2\mathbf{R}_1\mathbf{P}_0 + \mathbf{R}_4\mathbf{R}_3\mathbf{R}_2\mathbf{t}_1 + \mathbf{R}_4\mathbf{t}_3 + \mathbf{R}_4\mathbf{R}_3\mathbf{t}_2 + \mathbf{t}_4$$

Similar to the math we need to transform from the coordinate system of one camera to another.

Homogeneous coordinates

Homogeneous coordinates

Here is a point in 3D:

$$\mathbf{P} = \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix}$$

Homogeneous coordinates

Here is a point in 3D:

$$\mathbf{P} = \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix}$$

What if we made it more complicated and used four numbers?

$$\mathbf{P}_h = \begin{bmatrix} sP_x \\ sP_y \\ sP_z \\ s \end{bmatrix}$$

Uhm, okay?

So this means that

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \\ 6 \\ 2 \end{bmatrix}$$

are the same point in homogeneous coordinates

Euclidean transformations again

Rotation $\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$, with columns \mathbf{r}_i , and translation \mathbf{t}

$$\mathbf{P}_1 = \mathbf{R}\mathbf{P}_0 + \mathbf{t}$$

Euclidean transformations again

Rotation $\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$, with columns \mathbf{r}_i , and translation \mathbf{t}

$$\mathbf{P}_1 = \mathbf{R}\mathbf{P}_0 + \mathbf{t} = \mathbf{r}_1 P_x + \mathbf{r}_2 P_y + \mathbf{r}_3 P_z + \mathbf{t}$$

Euclidean transformations again

Rotation $\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$, with columns \mathbf{r}_i , and translation \mathbf{t}

$$\mathbf{P}_1 = \mathbf{R}\mathbf{P}_0 + \mathbf{t} = \mathbf{r}_1 P_x + \mathbf{r}_2 P_y + \mathbf{r}_3 P_z + \mathbf{t}$$

$$\mathbf{P}_1 = \begin{bmatrix} \mathbf{r}_1 & \mathbf{r}_2 & \mathbf{r}_3 & \mathbf{t} \end{bmatrix} \begin{bmatrix} P_x \\ P_y \\ P_z \\ 1 \end{bmatrix}$$

Euclidean transformations again

Rotation $\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$, with columns \mathbf{r}_i , and translation \mathbf{t}

$$\mathbf{P}_1 = \mathbf{R}\mathbf{P}_0 + \mathbf{t} = \mathbf{r}_1 P_x + \mathbf{r}_2 P_y + \mathbf{r}_3 P_z + \mathbf{t}$$
$$\mathbf{P}_1 = \underbrace{[\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3 \ \mathbf{t}]}_{\text{Transform: } \tilde{\mathbf{T}}} \underbrace{\begin{bmatrix} P_x \\ P_y \\ P_z \\ 1 \end{bmatrix}}_{\text{Homogeneous: } \mathbf{P}_{0h}} = \tilde{\mathbf{T}}\mathbf{P}_{0h}$$

Homogeneous euclidean transformations

Fully homogeneous euclidean transformations become

$$\mathbf{P}_{1h} = \mathbf{T}\mathbf{P}_{0h}$$

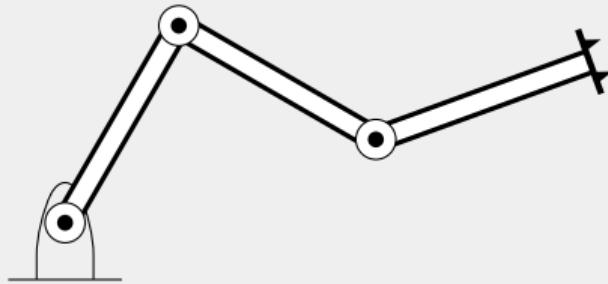
$$\begin{bmatrix} \mathbf{P}_1 \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{P}_0 \\ 1 \end{bmatrix}$$

Homogeneous euclidean transformations

The homogeneous transformation \mathbf{T} takes on the 4×4 form

$$\mathbf{T} = \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

Robot arm and homogeneous transformations



$$Q_{4h} = T_4 T_3 T_2 T_1 Q_{0h}$$

Not just easier; It is faster! and let's us do a lot of mathemagic.

**We can represent a rigid transformation
both as a 3×4 and as a 4×4 matrix.**

The *inhomogeneous* p corresponds to the
homogeneous p_h by

$$\mathbf{p}_h = \begin{bmatrix} \textcolor{red}{s} \mathbf{p} \\ \textcolor{red}{s} \end{bmatrix}$$

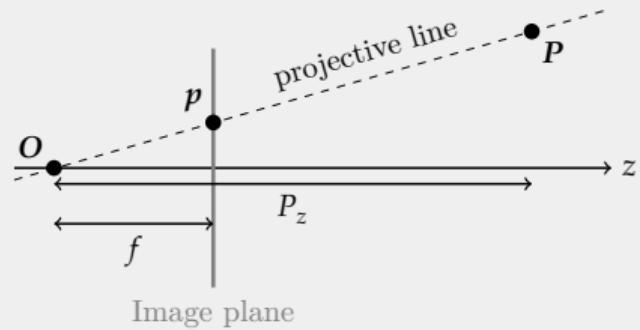
Questions? Short break

The projective transformation

Assume $f = 1$:

$$p_x = \frac{P_x}{P_z}, \quad p_y = \frac{P_y}{P_z}$$

$$\mathbf{P} = \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix} = \begin{bmatrix} \textcolor{red}{s}p_x \\ \textcolor{red}{s}p_y \\ \textcolor{red}{s} \end{bmatrix} = \mathbf{p}_h$$



Projective transformation is like assuming point in 3D is a 2D homogeneous point.

**There are many different notations for
homogeneous coordinates**

$$\mathbf{q} = \begin{bmatrix} s\mathbf{p} \\ s \end{bmatrix}$$

Lines in homogenous coordinates

Consider the line given implicitly by

$$0 = ax + by + c$$

We can write this in homogeneous coordinates

$$\begin{aligned} &= \begin{bmatrix} a \\ b \\ c \end{bmatrix}^T \begin{bmatrix} sx \\ sy \\ s \end{bmatrix} \\ &= \mathbf{l}^T \mathbf{p}_h \end{aligned}$$

Lines in homogenous coordinates

If $a^2 + b^2 = 1$ and the scale of the homogeneous point is 1 then,

$$d = \begin{bmatrix} a \\ b \\ c \end{bmatrix}^T \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = l^T p_h$$

d is the **signed distance** from the point to the line.

The homogeneous coordinate system — summary

The additional imaginary scale s , or alternatively dimension w

$$\mathbf{u} = s \begin{bmatrix} \mathbf{v} \\ 1 \end{bmatrix} = \begin{bmatrix} s\mathbf{v} \\ s \end{bmatrix} = \begin{bmatrix} \mathbf{v}' \\ w \end{bmatrix}$$

- Dimensionality is $N + 1$
- Points have a scale $s \neq 0$
- Directions have $w = 0$
- Many-to-one correspondence: $\mathbf{u} \in \mathbb{R}^{N+1}$ and $\mathbf{v} \in \mathbb{R}^N$

Getting *inhomogeneous* coordinates back

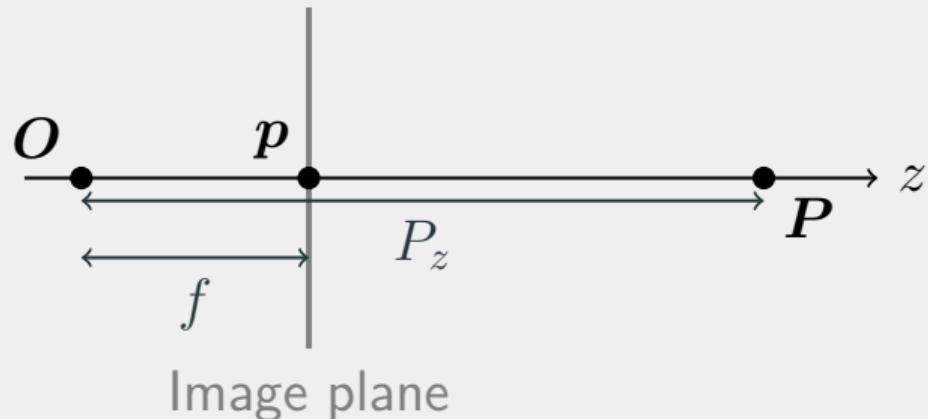
$$\mathbf{v} = \Pi(\mathbf{u}) = \Pi\left(\begin{bmatrix} \mathbf{v}' \\ s \end{bmatrix}\right) = \mathbf{v}'/s$$

Trivial inverse

Pinhole camera model

Principal point

Let P be a point exactly in the viewing direction of the camera.



p is called the **principal point**.

Where is p in the image on the right?

Principal point

Where is $(0, 0)$ in the image?

- After projective transformation

Principal point

Where is $(0, 0)$ in the image?

- After projective transformation
 - At the principal point

Principal point

Where is $(0, 0)$ in the image?

- After projective transformation
 - At the principal point
- Pixel coordinates
 - Upper left corner

We introduce δ_x and δ_y to translate $(0, 0)$ from the principal point to the upper left corner.

It is typically around half of the resolution of the camera.

Principal point

Projection is now

$$p_x = \frac{f}{P_z} P_x + \delta_x$$

$$p_y = \frac{f}{P_z} P_y + \delta_y$$

Can we write all of this using homogeneous coordinates?

Principal point

Yes we can!

$$\mathbf{p}_h = \underbrace{\begin{bmatrix} f & 0 & \delta_x \\ 0 & f & \delta_y \\ 0 & 0 & 1 \end{bmatrix}}_K \mathbf{P}$$

And it even looks nice! This is called the camera matrix.

It contains intrinsic camera parameters.

Extrinsics

- So far we have assumed the camera is at the origin $(0, 0, 0)$
- Does not generalize well to multiple cameras
- Solution:
 - Introduce a canonical “world” coordinate system
 - Transform points from world to camera before projecting

Extrinsics

- We parametrize this transformation with a
 - \mathbf{R} (rotation)
 - \mathbf{t} (translation)
- These are the **extrinsics**
- Transforming to the reference frame of the camera:

$$\begin{aligned}\mathbf{P}_{cam} &= \mathbf{RP} + \mathbf{t} \\ &= [\mathbf{R} \ \mathbf{t}] \begin{bmatrix} \mathbf{P} \\ 1 \end{bmatrix}\end{aligned}$$

Projection matrix

Projecting a single point in 3D to the camera

$$\mathbf{p}_h = \mathbf{K} \mathbf{P}_{cam}$$

Projection matrix

Projecting a single point in 3D to the camera

$$\begin{aligned} \mathbf{p}_h &= \mathbf{K} \mathbf{P}_{cam} \\ &= \underbrace{\mathbf{K} [\mathbf{R} \quad \mathbf{t}]}_{\mathcal{P}} \mathbf{P}_h \end{aligned}$$

Projection matrix

Projecting a single point in 3D to the camera

$$\begin{aligned} \mathbf{p}_h &= \mathbf{K} \mathbf{P}_{cam} \\ &= \underbrace{\mathbf{K} [\mathbf{R} \quad \mathbf{t}]}_{\mathcal{P}} \mathbf{P}_h = \end{aligned}$$

$$\begin{bmatrix} sp_x \\ sp_y \\ s \end{bmatrix} = \begin{bmatrix} f & 0 & \delta_x \\ 0 & f & \delta_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \end{bmatrix} \begin{bmatrix} P_x \\ P_y \\ P_z \\ 1 \end{bmatrix}$$

Wrapping up

- The translation is not the position of the camera
- What happens if the last coordinate of P_h is not 1?

Wrapping up

- The translation is not the position of the camera
- What happens if the last coordinate of P_h is not 1?
 - We get a scaled version of P_{cam} but, it is along the same line, so it projects to the same point.

Projection matrix:

$$q = \mathcal{P}P_h = K [R \ t] P_h$$

The matrix \mathcal{P} is known as the
projection matrix
(don't call it the camera matrix)

Exercise information

- Use Python interactively
 - Jupyter notebook
 - VS Code
 - Spyder
 - etc.
- Makes it easier to debug

Exercise information: Storing points on the computer

- Storing multiple one-dimensional vectors happens frequently
- Matrices are ideal for this
- We always operate on column vectors, so these matrices should be $3 \times n$ for many 3D vector (for example)
- Matrix multiplication lets you project many points at once
 - $\text{ph} = P \cdot \text{dot}(Ph)$ or even shorter
 - $\text{ph} = P @ Ph$

Comment about exercises

- NumPy is your friend!
- If you need for-loops, you're probably not doing it the easy way.
 - No exercises today need for-loops (except the provided function)
- Converting from homogeneous coordinates to regular
 - $p = ph[:-1]/ph[-1]$
- Ask the TAs 😊

Learning objectives

After this lecture you should be able to:

- explain homogeneous coordinates
- convert to and from homogeneous coordinates
- perform relevant coordinate transformations
- explain the pinhole camera model

Exercise time!

Camera model and homographies

Camera parameters; lens distortion; homography estimation

Morten R. Hannemose, mohan@dtu.dk

February 9, 2024

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Quiz 1 - Results

Correct		
Question 1	Line distance	67%
Question 2	Projection	79%

Add your name this week to join a leaderboard.

Learning objectives

After this lecture you should be able to:

- explain the phenomenological origin of lens distortion
- compensate for lens distortion in images
- explain the phenomenological origin of the pinhole camera parameters
- derive the homography matrix
- explain the use of homographies in computer vision
- perform the linear estimation of a homography matrix
- use singular value decomposition (SVD) to solve linear systems

Presentation topics

Pinhole camera

Intrinsics

Lens distortion

Radial distortion

Homographies

Application: image stitching (panorama)

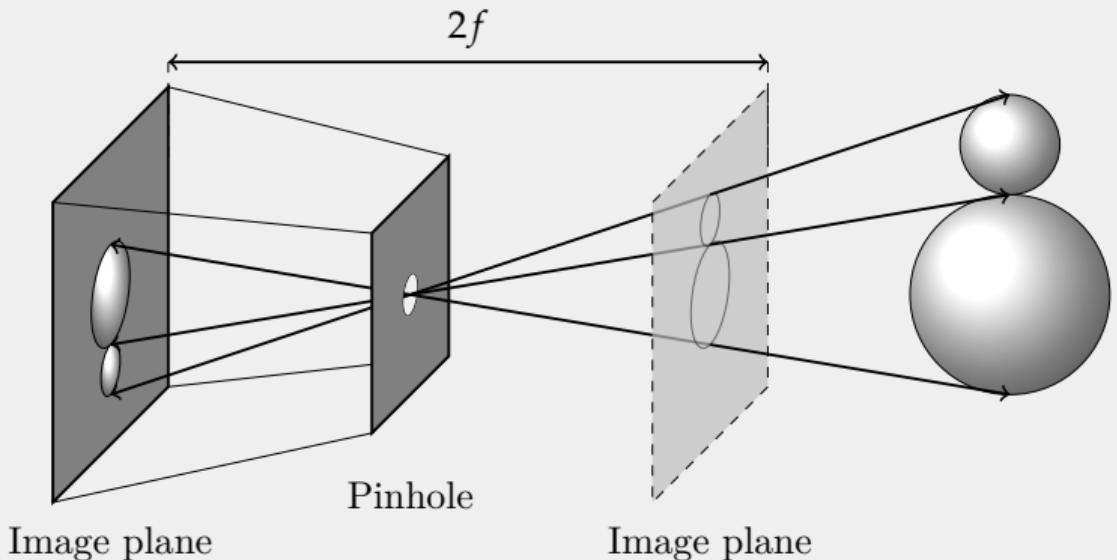
Homography estimation

Homogeneous least-squares solution

Normalization of points

Pinhole camera

Recap last week



$$\mathbf{K} = \begin{bmatrix} f & 0 & \delta_x \\ 0 & f & \delta_y \\ 0 & 0 & 1 \end{bmatrix}$$

Camera intrinsics

- Focal length f
- Principal point $\delta x, \delta y$
- Skew parameters α and β

Focal length f

- The focal length is the distance from the pinhole to the image plane *in pixels*.
- The camera matrix makes projected coordinates equal pixels
 - δ_x and δ_y translates to make $(0, 0)$ correct
 - f performs the scaling to make the projected point into pixel coordinates

- The focal length describes the level of zoom / field of view
 - Longer focal length = more zoomed in.

Focal length f

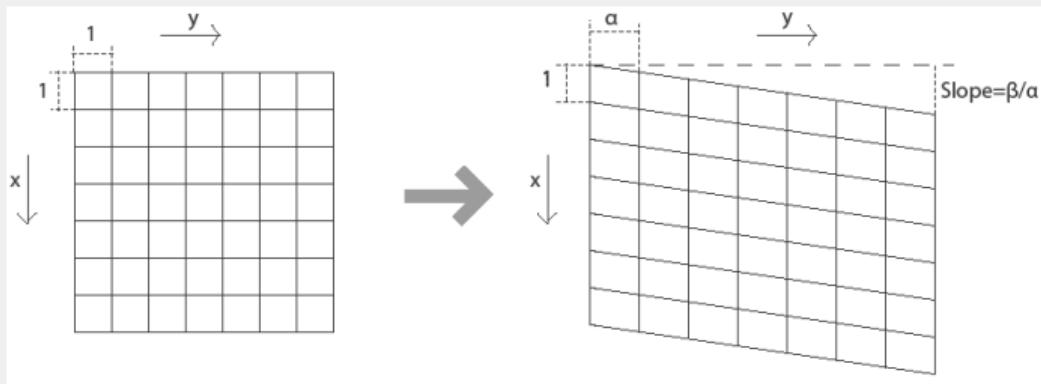
- How is the image affected visually by changing the focal length?
- Moving camera closer while zooming out to make one object have a constant size is a popular effect
- Real-life example
- Simple example
- This is called a dolly zoom / the vertigo effect

Why $\begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix}$ **and not** $\begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & f \end{bmatrix}$?

Skew parameters α and β

Non-square pixels (α),
and non-rectangular
pixels (β)

$$\boldsymbol{K} = \begin{bmatrix} f & \beta f & \delta_x \\ 0 & \alpha f & \delta_y \\ 0 & 0 & 1 \end{bmatrix}$$



Orthographic projection

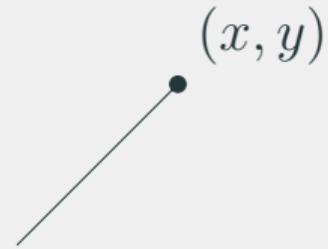
If all projection lines are parallel, then magnification $m = 1$

$$\mathbf{p} = \begin{bmatrix} P_x \\ P_y \end{bmatrix}$$

The pinhole model goes towards orthographic projection as the focal length goes to ∞

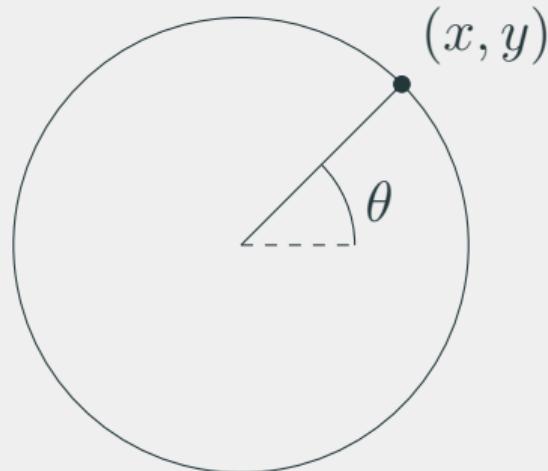
Degrees of Freedom

A point in 2D



Degrees of Freedom

A point in 2D



has one degree of freedom when it lies on a circle.

Projection matrix: Degrees of freedom

$$\mathcal{P} = \underbrace{\begin{bmatrix} f & \beta f & \delta_x \\ 0 & \alpha f & \delta_y \\ 0 & 0 & 1 \end{bmatrix}}_K \underbrace{\begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix}}_{[\boldsymbol{R} \quad \boldsymbol{t}]}$$

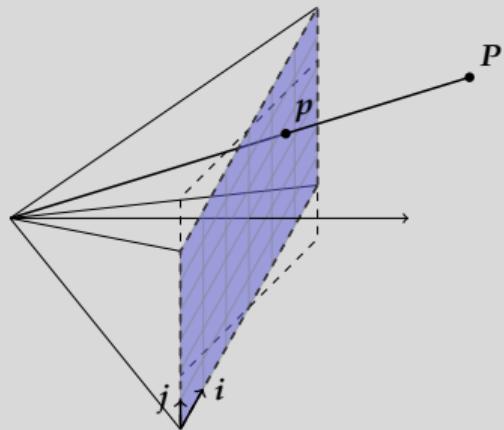
Projection matrix: Degrees of freedom

$$\mathcal{P} = \underbrace{\begin{bmatrix} f & \beta f & \delta_x \\ 0 & \alpha f & \delta_y \\ 0 & 0 & 1 \end{bmatrix}}_K \underbrace{\begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix}}_{[\mathbf{R} \quad \mathbf{t}]}$$

\mathcal{P} is $3 \times 4 = 12$ numbers.

$5 + 3 + 3 = 11$ degrees of freedom + 1 (scale) = 12

Which physical properties do I need to derive the intrinsics: f , δx , δy , α , β ?



**Speculate; Given an unknown system,
how would you estimate all parameters in
 \mathcal{P} ?**

Lens distortion

Help! My camera is not a pinhole camera?



Welcome to the real world!



@OPPO

Real lenses

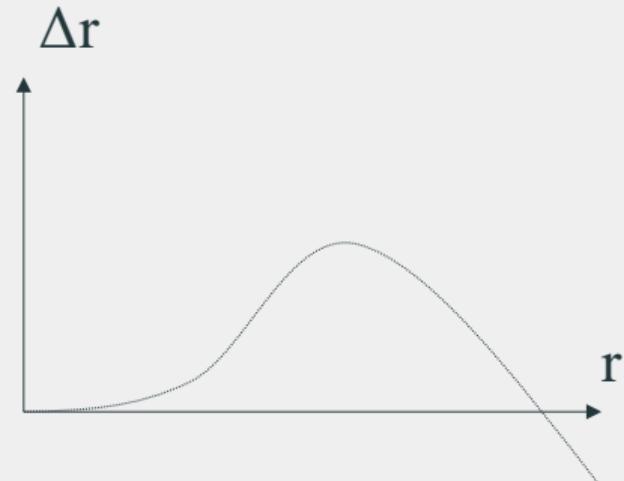
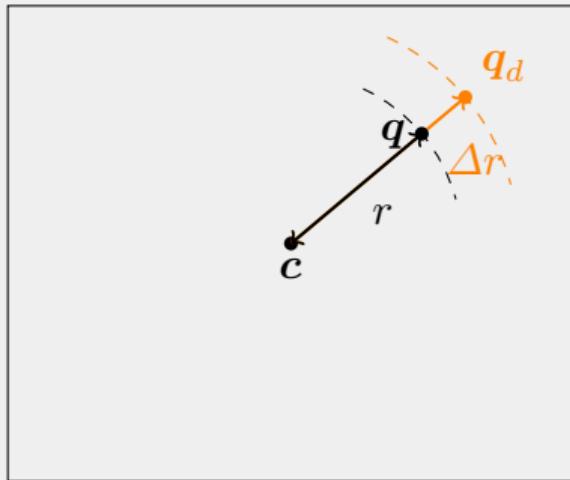
Lenses are:

- complicated by many effects, intentional or non-intentional
- imperfect with, for example, defects
- different for each camera/lens

Almost all lenses have lens distortion!

Radial lens distortion

We distort points by a polynomial in the distance from the principal point (the radius r)



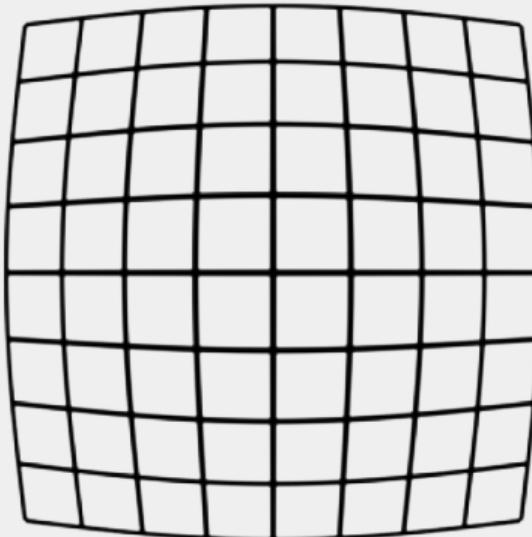
Radial distortion equations

$$\begin{bmatrix} x_d \\ y_d \end{bmatrix} = \text{dist} \left(\begin{bmatrix} x \\ y \end{bmatrix} \right) = \begin{bmatrix} x \\ y \end{bmatrix} \cdot \left(1 + \Delta r(\sqrt{x^2 + y^2}) \right)$$
$$\Delta r(r) = k_3 r^2 + k_5 r^4 + k_7 r^6 \dots$$

Barrel distortion

For example:

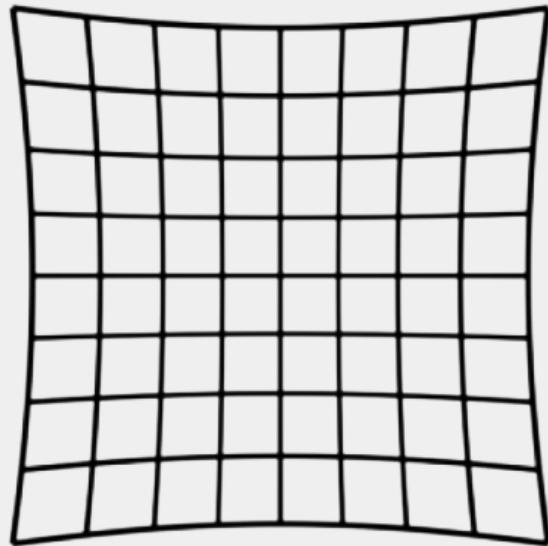
$$\Delta r(r) = -0.5r^2$$



Pincushion distortion

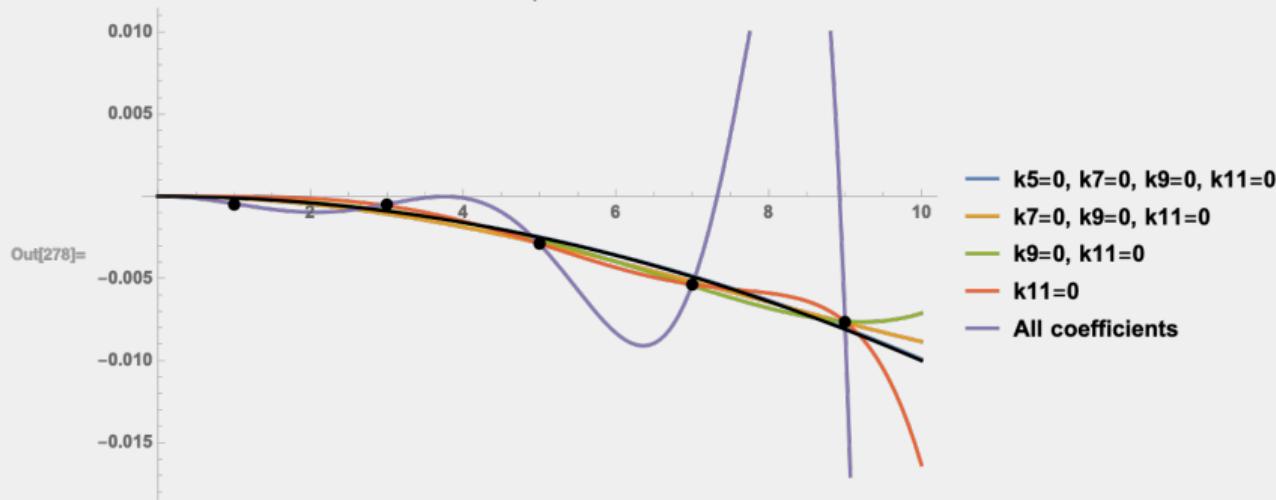
For example:

$$\Delta r(r) = 0.5r^2$$



Radial lens distortion in practice

Overfitting / extrapolation errors



$$\Delta r(r) = -0.1r^2 + \delta_{error}$$

Almost all lenses have lens distortion!

Projection under radial distortion

Projection without distortion:

$$p = K \underbrace{[R|t]Q}_q$$

Projection under radial distortion

Projection without distortion:

$$p = K \underbrace{[R|t]Q}_q = K\Pi^{-1} \left(\Pi \left(\underbrace{[R|t]Q}_q \right) \right)$$

Projection under radial distortion

Projection without distortion:

$$p = \underbrace{\mathbf{K} [\mathbf{R}|t]\mathbf{Q}}_q = \mathbf{K} \Pi^{-1} \left(\Pi \left(\underbrace{[\mathbf{R}|t]\mathbf{Q}}_q \right) \right)$$

Projection with distortion:

$$\mathbf{p}_d = \underbrace{\mathbf{K} \Pi^{-1} \left(\text{dist} \left(\Pi \left([\mathbf{R}|t]\mathbf{Q} \right) \right) \right)}_{\mathbf{q}_d}$$

Π converts from homogeneous to inhomogeneous coordinates.
 $\text{dist}(\cdot)$ operates on 2D points.

Radial distortion equations [recap]

$$\begin{bmatrix} x_d \\ y_d \end{bmatrix} = \text{dist} \left(\begin{bmatrix} x \\ y \end{bmatrix} \right) = \begin{bmatrix} x \\ y \end{bmatrix} \cdot \left(1 + \Delta r(\sqrt{x^2 + y^2}) \right)$$
$$\Delta r(r) = k_3 r^2 + k_5 r^4 + k_7 r^6 \dots$$

Comparison to lecture notes

- This follows the math in OpenCV unlike the Lecture Notes
- Pros:
 - Projection can be done without splitting the camera matrix in two
 - Distortion parameters still work if we resize the image
- Cons:
 - Undistorting an image requires knowing the focal length

Undistortion

- You are given a distorted image, the camera matrix and the distortion coefficients
- Q1: Can you undistort points observed in the distorted image?
- Q2: Can you undistort the image?
- Discuss with the person next to you

Undistorting points (answer to Q1)

- The math given so far describes how to distort points.
- Undistorting points requires you to invert $\text{dist}(\cdot)$
- There is no analytical solution to this, so it must be done iteratively 😞

Undistorting an image (answer to Q2)

- Having the RGB-value of every (integer) pixel location (p) in the undistorted image is the same as having the image.
- However, we only know the RGB-value at every location (p_d) in the distorted image

Undistorting an image (answer to Q2)

- Having the RGB-value of every (integer) pixel location (p) in the undistorted image is the same as having the image.
- However, we only know the RGB-value at every location (p_d) in the distorted image
- Use $\text{dist}(\cdot)$ and K to find p_d from p
 - $p_d = K\Pi^{-1} \left(\text{dist} \left(\Pi \left(K^{-1}p \right) \right) \right)$
- Use bilinear interpolation to compute the RGB-value in the distorted image at the (non-integer) point p_d .
- Not necessary to invert the distortion function!

Questions? Short break

Homographies

Practical example

- No volunteers needed 😊

Practical example

- No volunteers needed 😊
- Knowing that a point seen with a camera lies on a plane lets us know exactly where on the plane it is.

The homography

Consider the following plane:

$$\begin{aligned} \mathbf{P} &= a\mathbf{A} + b\mathbf{B} + \mathbf{C}, \\ &= [\mathbf{A} \ \mathbf{B} \ \mathbf{C}] \begin{bmatrix} a \\ b \\ 1 \end{bmatrix}. \end{aligned}$$

\mathbf{C} is a point on the plane, and \mathbf{A} and \mathbf{B} are vectors that span the plane.

The homography

Projections of points on a plane:

$$\begin{aligned} \mathbf{q} &= \mathcal{P} \mathbf{P}_h, \\ &= \mathcal{P} \begin{bmatrix} A & B & C \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \\ 1 \end{bmatrix}, \end{aligned}$$

The homography

Projections of points on a plane:

$$\begin{aligned} \mathbf{q} &= \mathcal{P} \mathbf{P}_h, \\ &= \underbrace{\mathcal{P} \begin{bmatrix} A & B & C \\ 0 & 0 & 1 \end{bmatrix}}_H \begin{bmatrix} a \\ b \\ 1 \end{bmatrix}, \end{aligned}$$

\mathbf{H} is a 3×3 matrix called a **homography**.

The homography



$$\text{H} \curvearrowright$$



Homographies between images

Two cameras photographing the same plane let's us establish pixel correspondence:

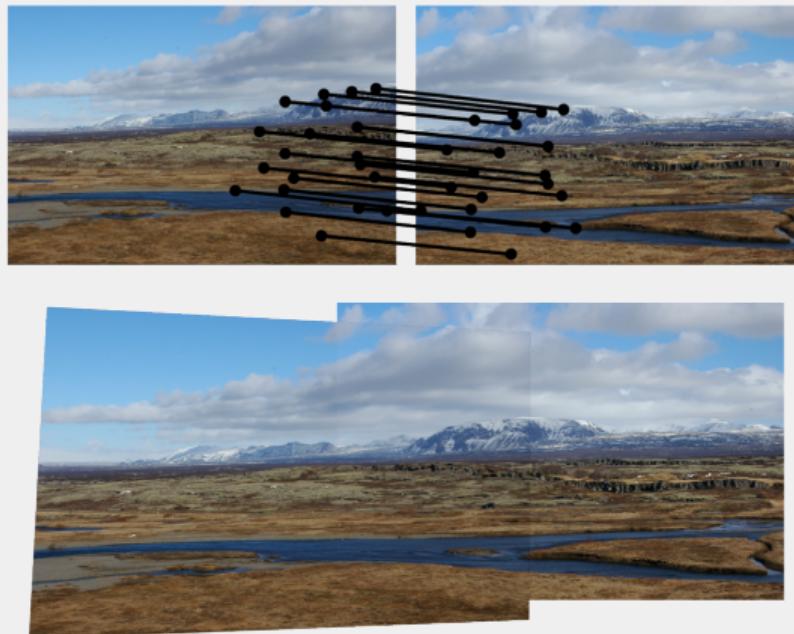
$$\mathbf{q}_1 = \mathbf{H}_1 \mathbf{Q}, \quad \mathbf{q}_2 = \mathbf{H}_2 \mathbf{Q}$$

$$\begin{aligned}\mathbf{q}_1 &= \mathbf{H}_1 \mathbf{H}_2^{-1} \mathbf{q}_2 \\ &= \mathbf{H} \mathbf{q}_2\end{aligned}$$

Two view geometry without baseline

If two cameras are at the same point we can use a homography to relate their pixels

Image stitching



You get to implement this in week 9.

Homography multiplication

If you multiply a homography H by a scalar, how does that change the mapping it describes?

Homography multiplication

If you multiply a homography H by a scalar, how does that change the mapping it describes?

It doesn't!

$$q_1 = Hq_2$$

$$sq_1 = (sH)q_2$$

q_1 and sq_1 are the same point.

Homography degrees of freedom

- 3×3 matrix has 9 numbers
- Scale is arbitrary
 - A homography has 8 degrees of freedom
- We need 4 pairs of points to estimate a homography
 - Each point pair imposes two constraints (x and y)

Any two images of the same plane are
related by a homography

Homography estimation

Homography estimation

Consider a number of points q_{1i} and q_{2i} . The homography relates them as follows:

$$q_{1i} = Hq_{2i}$$

Homography rewritten i

The cross product of two parallel vectors is the zero vector:

$$\begin{aligned}\mathbf{q}_{1i} &= \mathbf{H}\mathbf{q}_{2i} \Rightarrow \\ \mathbf{0} &= \mathbf{q}_{1i} \times \mathbf{H}\mathbf{q}_{2i},\end{aligned}$$

we can isolate \mathbf{H} in this expression . . .

Homography rewritten ii

$$\mathbf{0} = \mathbf{q}_{1i} \times \mathbf{H} \mathbf{q}_{2i},$$

$$= \mathbf{B}^{(i)} \text{flatten}(\mathbf{H}^T),$$

$$\mathbf{B}^{(i)} = \begin{bmatrix} 0 & -x_{2i} & x_{2i}y_{1i} & 0 & -y_{2i} & y_{2i}y_{1i} & 0 & -1 & y_{1i} \\ x_{2i} & 0 & -x_{2i}x_{1i} & y_{2i} & 0 & -y_{2i}x_{1i} & 1 & 0 & -x_{1i} \\ -x_{2i}y_{1i} & x_{2i}x_{1i} & 0 & -y_{2i}y_{1i} & y_{2i}x_{1i} & 0 & -y_{1i} & x_{1i} & 0 \end{bmatrix},$$

$$\text{flatten}(\mathbf{H}^T) = [H_{11} \ H_{21} \ H_{31} \ H_{12} \ H_{22} \ H_{32} \ H_{13} \ H_{23} \ H_{33}]^T$$

Defining $B^{(i)}$

The $B^{(i)}$ matrices can be implemented with the Kronecker product
np.kron

$$\begin{aligned}B^{(i)} &= \mathbf{q}_{2i}^T \otimes [\mathbf{q}_{1i}]_{\times} \\&= \mathbf{q}_{2i}^T \otimes \begin{bmatrix} 0 & -1 & y_{1i} \\ 1 & 0 & -x_{1i} \\ -y_{1i} & x_{1i} & 0 \end{bmatrix}\end{aligned}$$

Homography matrix solution

Given n corresponding points, we stack their respective $\mathbf{B}^{(i)}$ matrices into a tall matrix \mathbf{B} :

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}^{(1)} \\ \mathbf{B}^{(2)} \\ \vdots \\ \mathbf{B}^{(n)} \end{bmatrix}$$
$$\mathbf{0} = \mathbf{B} \text{ flatten}(\mathbf{H}^T)$$

How do we solve this linear system of equations?

Homogeneous least-squares problem

- Problem of type: $\mathbf{A}\mathbf{x} = \mathbf{0}$
- When noise is present, cannot be satisfied exactly
- Equivalent to $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x}\|_2^2$
- Trivial solution $\mathbf{x} = \mathbf{0}$ (not useful)
- \mathbf{x} can be arbitrarily scaled
 - We impose $\|\mathbf{x}\|_2 = 1$

Homogeneous least-squares solution

Rewrite problem

$$\begin{aligned}\|Ax\|_2^2 &= \\ (Ax)^T Ax &= \\ x^T A^T Ax\end{aligned}$$

Assume x is an **eigenvector** of the square matrix $A^T A$.

Homogeneous least-squares solution

Which eigenvector?

Eigen-decomposition of $\mathbf{A}^\top \mathbf{A}$:

$$\mathbf{A}^\top \mathbf{A} = \mathbf{V} \Lambda \mathbf{V}^\top,$$
$$\mathbf{v}_i^\top \mathbf{A}^\top \mathbf{A} \mathbf{v}_i = \underbrace{\mathbf{v}_i^\top \mathbf{V}}_{e_i^\top} \Lambda \underbrace{\mathbf{V}^\top \mathbf{v}_i}_{e_i} = \lambda_i,$$

where e_i is the vector with 0 everywhere except for the i^{th} index.

To minimize λ_i we choose \mathbf{v}_i (and \mathbf{x}) as the eigenvector with the smallest eigenvalue.

Least-squares solution with SVD

We can use the SVD (`np.linalg.svd`)

$$\begin{aligned} \mathbf{A}^\top \mathbf{A} &= (\mathbf{U} \Sigma \mathbf{V}^\top)^\top \mathbf{U} \Sigma \mathbf{V}^\top \\ &= \mathbf{V} \Sigma^\top \mathbf{U}^\top \mathbf{U} \Sigma \mathbf{V}^\top \\ &= \mathbf{V} \underbrace{\Sigma^\top \Sigma}_{\Lambda} \mathbf{V}^\top \end{aligned}$$

\mathbf{V} are the eigenvectors of $\mathbf{A}^\top \mathbf{A}$, and the minimum solution has the **smallest singular value**.

Least-squares solution with SVD

- When we use the SVD, we avoid computing $A^T A$
- The least-squares solution is
 - the right singular vector of A corresponding the smallest singular value
 - the eigenvector of $A^T A$ corresponding to the smallest eigenvalue

The solution of $Ax = 0$ where $\|x\|_2 = 1$ is found using SVD.

$x = v$ where the right singular vector v corresponds to the smallest singular value.

Normalization of points

Numerical instability

Experience has shown that most of these algorithms can be numerically unstable.

What can we do? We can normalize!

Point normalization

Consider the collection of n points \mathbf{p}_i that has mean and standard deviation as follows

$$\text{mean}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n) = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} = \boldsymbol{\mu} \text{ and}$$
$$\text{std}(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n) = \begin{bmatrix} \sigma_x \\ \sigma_y \end{bmatrix} = \boldsymbol{\sigma}$$

Point normalization

We would like to normalize these points so they have mean **0** and standard deviation **1**.

Let us call the normalized points \tilde{p}_i .

$$\text{mean}(\tilde{p}_1, \tilde{p}_2, \dots, \tilde{p}_n) = 0 \text{ and}$$

$$\text{std}(\tilde{p}_1, \tilde{p}_2, \dots, \tilde{p}_n) = 1$$

Point normalization

The relation between p_i and \tilde{p}_i is then

$$\begin{aligned}\tilde{p}_i &= \frac{p_i - \mu}{\sigma} \Leftrightarrow \\ p_i &= \sigma \tilde{p}_i + \mu\end{aligned}$$

Can we write this with homogeneous coordinates with some transformation T ?

$$\tilde{q}_i = T q_i$$

Point normalization in homogeneous coordinates

$$\tilde{\mathbf{q}}_i = \underbrace{\begin{bmatrix} 1/\sigma_x & 0 & -\mu_x/\sigma_x \\ 0 & 1/\sigma_y & -\mu_y/\sigma_y \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{T}} \mathbf{q}_i$$

$$\mathbf{q}_i = \underbrace{\begin{bmatrix} \sigma_x & 0 & \mu_x \\ 0 & \sigma_y & \mu_y \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{T}^{-1}} \tilde{\mathbf{q}}_i$$

It is easiest to implement \mathbf{T}^{-1} and invert it to obtain \mathbf{T} .

Point normalization in linear algorithms

Now we can estimate the homography using \tilde{p}_{1i} and \tilde{p}_{2i} and obtain the homography $\widetilde{\mathbf{H}}$.

How do we get the homography \mathbf{H} that operates on the original points?

$$\begin{aligned}\tilde{\mathbf{q}}_{1i} &= \widetilde{\mathbf{H}}\tilde{\mathbf{q}}_{2i} = \\ \mathbf{T}_1 \mathbf{q}_{1i} &= \widetilde{\mathbf{H}} \mathbf{T}_2 \mathbf{q}_{2i} \\ \mathbf{q}_{1i} &= \underbrace{\mathbf{T}_1^{-1} \widetilde{\mathbf{H}} \mathbf{T}_2}_\mathbf{H} \mathbf{q}_{2i}\end{aligned}$$

Exercise tips

- You make mistakes and you learn from them.
- Follow these good practice tips:
- `im = cv2.imread('im.jpg')[:, :, ::-1]`

Exercise tips

- You make mistakes and you learn from them.
- Follow these good practice tips:
- `im = cv2.imread('im.jpg')[:, :, ::-1]`
- `im = im.astype(float)/255`

Learning objectives

After this lecture you should be able to:

- explain the phenomenological origin of lens distortion
- compensate for lens distortion in images
- explain the phenomenological origin of the pinhole camera parameters
- derive the homography matrix
- explain the use of homographies in computer vision
- perform the linear estimation of a homography matrix
- use singular value decomposition (SVD) to solve linear systems

Exercise time!

Stereo view geometry

Epipolar geometry, triangulation

Morten R. Hannemose, mohan@dtu.dk

February 15, 2024

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- Derive and explain the epipolar line in computer vision
- Derive and apply the fundamental matrix in computer vision
- Derive and apply the essential matrix in computer vision
- Implement the linear algorithm for triangulation
- Explain the pros and cons of using a linear algorithm

Presentation topics

Projection of points, lines and planes

Epipolar planes and lines

Essential and fundamental matrices

Triangulation

Problems with linear algorithm

Projection of points, lines and planes

Point projections

Any 3D point P projects to a 2D point p :

$$P \xrightarrow{\text{projection}} p.$$

Line projections

Any 3D line \mathbf{L} projects to a 2D line \mathbf{l} , except if $\mathbf{L} \parallel \mathbf{L}_{\mathcal{P}}$:

$$\mathbf{L}_0 \xrightarrow{\text{projection}} \mathbf{l}_0 ,$$

$$\mathbf{L}_1 \xrightarrow{\text{projection}} \mathbf{p}_1 .$$

Plane projections

Any 3D plane P projects to
the image plane,

$$P \xrightarrow{\text{projection}} p,$$

except if $P \parallel L_\varphi$
Where then?

Epipolar planes and lines

Epipolar lines

- Throwback to the demonstration from the first week
- A point seen in one camera must lie along a specific 3D line
- A point seen in one camera must lie along a specific 2D line in the other camera
 - This is the epipolar line!

Stereo view

Assume stereo view (two cameras)

Epipolar planes

The camera centers (q_1 and q_2) and the 3D point Q form an

Epipolar planes

Which line in space is always part of the epipolar planes?

Epipolar lines

The epipolar lines intersect the epipolar plane and images!

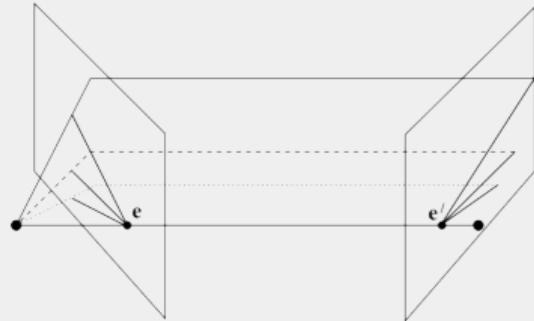
Epipolar lines

The relation between the line $q_i Q$ and the epipolar lines?

Epipolar lines

All epipolar lines intersect in the epipoles.

Epipoles



Where are the epipoles?

**Stereo view is an epipolar geometry:
Each Q in 3D has an epipolar plane.**

**Pixels correspond if and only if
they lie in the same epipolar plane**

Essential and fundamental matrices

Setting the scene

Let q refer to a 2D point in pixels and p refer to the same point in 3D in the reference frame of the camera.

These are related as follows:

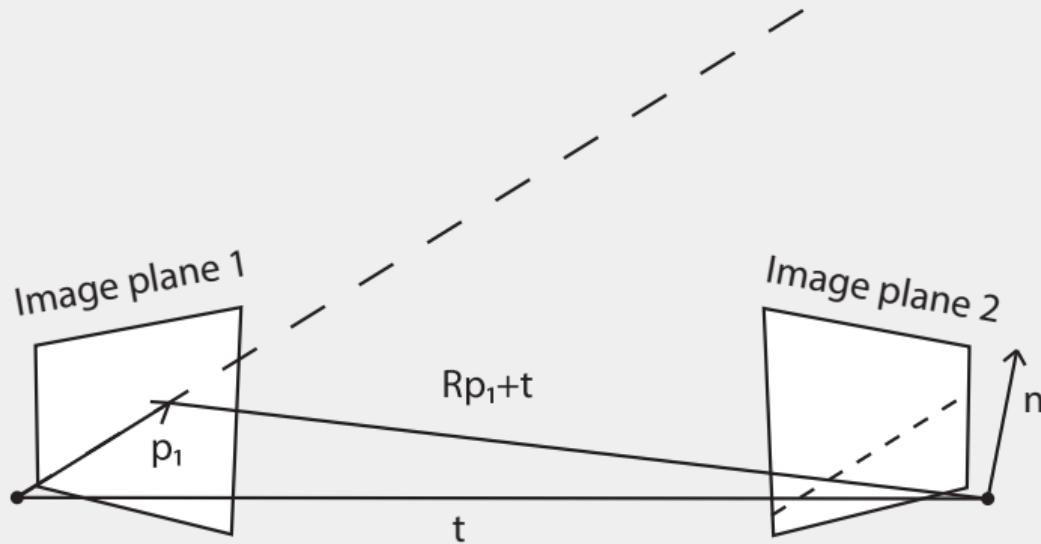
$$q = Kp$$

$$p = K^{-1}q$$

R, t maps from the reference frame of camera one to the reference frame of camera two (their relative transformation)

The essential matrix

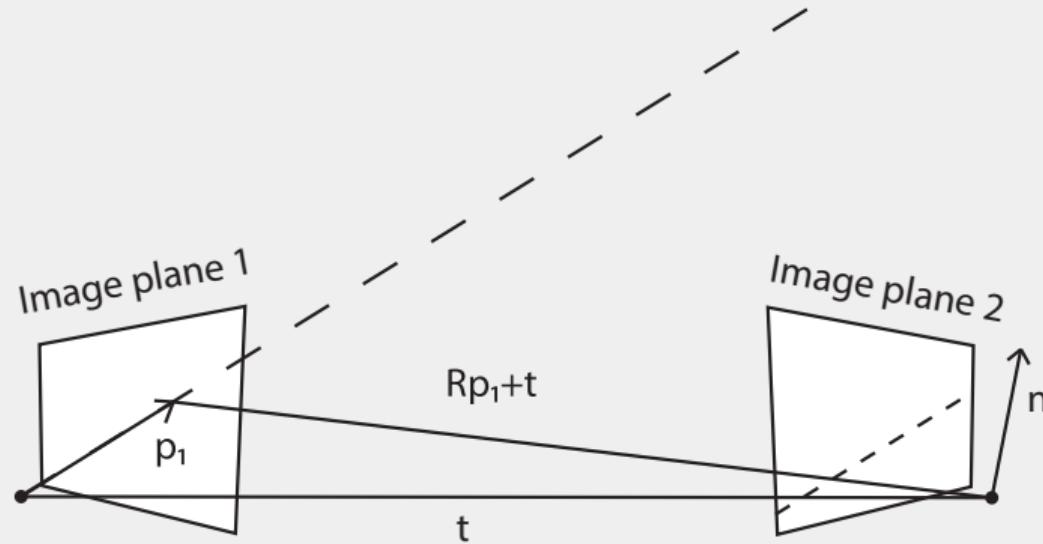
Consider the epipolar plane given by p_1



Vectors in the figure are in the reference frame of camera 2

The essential matrix

Relative to camera 2, what is the normal of the epipolar plane?



The normal is orthogonal to t and $Rp_1 + t$

The essential matrix

The normal is then

$$\begin{aligned}\mathbf{n} &= \mathbf{t} \times (\mathbf{R}\mathbf{p}_1 + \mathbf{t}) \\ &= \mathbf{t} \times (\mathbf{R}\mathbf{p}_1) \\ &= \underbrace{[\mathbf{t}]_{\times}}_{\mathbf{E}} \mathbf{R} \mathbf{p}_1.\end{aligned}$$

\mathbf{E} is called the **essential matrix**!

The essential matrix

Dot product of orthogonal vectors are zero.

For the corresponding p_2 in the second camera

$$\begin{aligned} 0 &= p_2^T n \\ &= p_2^T E p_1 \end{aligned}$$

The essential matrix imposes a constraint on corresponding points.

How to interpret this?

- $p = K^{-1}q$.
- Are p_1 and p_2 in homogeneous or inhomogeneous coordinates?

How to interpret this?

- $p = K^{-1}q$.
- Are p_1 and p_2 in homogeneous or inhomogeneous coordinates?
 - Yes!...
- There are two interpretations:

How to interpret this?

- $p = K^{-1}q$.
- Are p_1 and p_2 in homogeneous or inhomogeneous coordinates?
 - Yes!...
- There are two interpretations:
 - p_1 and p_2 are 3D points and n is a vector in 3D. We use this in our derivations.
 - p_1 and p_2 are 2D points and $n = Ep_1$ is the epipolar line, both are in homogeneous coordinates.
- 

The fundamental matrix

Recall that

$$\mathbf{p} = \mathbf{K}^{-1}\mathbf{q}.$$

Then

$$\begin{aligned}\mathbf{p}_2^\top \mathbf{E} \mathbf{p}_1 &= 0 \\ (\mathbf{K}_2^{-1}\mathbf{q}_2)^\top \mathbf{E} (\mathbf{K}_1^{-1}\mathbf{q}_1) &= 0 \\ \mathbf{q}_2^\top \underbrace{\mathbf{K}_2^{-\top} \mathbf{E} \mathbf{K}_1^{-1}}_F \mathbf{q}_1 &= 0\end{aligned}$$

where F is the **fundamental matrix**!

The essential and fundamental matrices form requirements for pixel correspondence:

$$p_2^T E p_1 = 0$$

$$q_2^T F q_1 = 0$$

Note on R , t

What if R , t is not given, but you only know the pose of each camera in world coordinates (R_1 , t_1 , R_2 , t_2)?

You can compute the relative transformation (which you will do in the exercise today)

Fundamental/essential matrix vs homography

What is the difference?

Fundamental/essential matrix vs homography

What is the difference?

- The fundamental and essential matrices yield epipolar lines:
 - $0 = \mathbf{p}_2^T \mathbf{E} \mathbf{p}_1$
 - $0 = \mathbf{q}_2^T \mathbf{F} \mathbf{q}_1$
 - Corresponding points must lie on the epipolar line, no matter what is in the image.
- The homography establishes one-to-one correspondences:
 - $\mathbf{q}_2 = \mathbf{H} \mathbf{q}_1$
 - Only valid for planes.

Degrees of freedom of F

- F has 9 numbers, and is scale invariant.
- $[t]_\times$ has rank 2, and thus F has as well.
 - In other terms: $\det(F) = 0$.
- You will estimate F in week 9.

Short break

Triangulation

Triangulation



We've seen the same point in many (known) cameras and want to find the point in 3D.

Because of noise, there is no exact solution.

Triangulation problem

Consider a projection matrix

$$\mathcal{P}_i = \begin{bmatrix} p_i^{(1)} \\ p_i^{(2)} \\ p_i^{(3)} \end{bmatrix}$$

where $p_i^{(1)}$ is the i^{th} row of \mathcal{P}_i .

Triangulation equations

Projection gives the pixels

$$\mathbf{q}_i = \begin{bmatrix} s_i x_i \\ s_i y_i \\ s_i \end{bmatrix} = \mathcal{P}_i \mathbf{Q} = \begin{bmatrix} p_i^{(1)} Q \\ p_i^{(2)} Q \\ p_i^{(3)} Q \end{bmatrix}$$

This is two constraints (x_i, y_i) in three equations.

Triangulation equations

As $s_i = \mathbf{p}_i^{(3)} \mathbf{Q}$, we have

$$(\mathbf{p}_i^{(3)} \mathbf{Q}) \begin{bmatrix} x_i \\ y_i \end{bmatrix} = \begin{bmatrix} \mathbf{p}_i^{(1)} \mathbf{Q} \\ \mathbf{p}_i^{(2)} \mathbf{Q} \end{bmatrix}$$

rearranged into

$$\begin{aligned} \mathbf{0} &= \begin{bmatrix} \mathbf{p}_i^{(3)} x_i - \mathbf{p}_i^{(1)} \\ \mathbf{p}_i^{(3)} y_i - \mathbf{p}_i^{(2)} \end{bmatrix} \mathbf{Q} \\ &= \mathbf{B}^{(i)} \mathbf{Q} \end{aligned}$$

Triangulation solution

Define \mathbf{B}

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}^{(1)} \\ \mathbf{B}^{(2)} \\ \vdots \\ \mathbf{B}^{(n)} \end{bmatrix} = \begin{bmatrix} \mathbf{p}_1^{(3)}x_1 - \mathbf{p}_1^{(1)} \\ \mathbf{p}_1^{(3)}y_1 - \mathbf{p}_1^{(2)} \\ \mathbf{p}_2^{(3)}x_2 - \mathbf{p}_2^{(1)} \\ \mathbf{p}_2^{(3)}y_2 - \mathbf{p}_2^{(2)} \\ \vdots \end{bmatrix}.$$

Use SVD to find $\arg \min_Q \|\mathbf{BQ}\|_2$, s.t. $\|\mathbf{Q}\|_2 = 1$.

Linear algorithm, hmm

- $\arg \min_Q \|BQ\|_2, \quad \text{s.t. } \|Q\|_2 = 1.$
- This is a linear algorithm used to solve the problem.
- Which error would we actually like to minimize?

Linear algorithm, hmm

- $\arg \min_Q \|BQ\|_2, \quad \text{s.t. } \|Q\|_2 = 1.$
- This is a linear algorithm used to solve the problem.
- Which error would we actually like to minimize?
 - Depends why we believe that there are errors
 - Errors in the observed pixel location

Triangulation: ideal error

- Let $[\tilde{x}_i \ \tilde{y}_i]$ and $[x_i \ y_i]$ refer to the observed and projected pixel coordinates, respectively.
- Which error would we like to minimize?

$$e_{\text{ideal}} = \sum_i \left\| \begin{bmatrix} \tilde{x}_i \\ \tilde{y}_i \end{bmatrix} - \begin{bmatrix} x_i \\ y_i \end{bmatrix} \right\|_2^2$$

Triangulation uh oh!

- Which error are actually we minimizing?

$$\begin{aligned} e_{\text{algebraic}} &= \sum_i \left\| \mathbf{B}^{(i)} \mathbf{Q} \right\|_2^2 \\ &= \sum_i \left\| \begin{bmatrix} \mathbf{p}_i^{(3)} \tilde{x}_i - \mathbf{p}_i^{(1)} \\ \mathbf{p}_i^{(3)} \tilde{x}_i - \mathbf{p}_i^{(2)} \end{bmatrix} \mathbf{Q} \right\|_2^2 \\ &= \sum_i \left\| \underbrace{\mathbf{p}_i^{(3)} \mathbf{Q}}_{s_i} \begin{bmatrix} \tilde{x}_i \\ \tilde{x}_i \end{bmatrix} - \begin{bmatrix} \mathbf{p}_i^{(1)} \\ \mathbf{p}_i^{(2)} \end{bmatrix} \mathbf{Q} \right\|_2^2 \\ &= \sum_i \left\| s_i \begin{bmatrix} \tilde{x}_i \\ \tilde{y}_i \end{bmatrix} - s_i \begin{bmatrix} x_i \\ y_i \end{bmatrix} \right\|_2^2 \end{aligned}$$

Error terms compared

- We can rewrite the terms slightly again

$$e_{\text{ideal}} = \sum_i (\tilde{x}_i - x_i)^2 + (\tilde{y}_i - y_i)^2$$

$$e_{\text{algebraic}} = \sum_i s_i^2 (\tilde{x}_i - x_i)^2 + s_i^2 (\tilde{y}_i - y_i)^2$$

- s_i is larger for cameras that are further from Q

Error terms compared

- We can rewrite the terms slightly again

$$e_{\text{ideal}} = \sum_i (\tilde{x}_i - x_i)^2 + (\tilde{y}_i - y_i)^2$$

$$e_{\text{algebraic}} = \sum_i s_i^2 (\tilde{x}_i - x_i)^2 + s_i^2 (\tilde{y}_i - y_i)^2$$

- s_i is larger for cameras that are further from Q
- Why don't we just divide by s_i ?

Error terms compared

- We can rewrite the terms slightly again

$$e_{\text{ideal}} = \sum_i (\tilde{x}_i - x_i)^2 + (\tilde{y}_i - y_i)^2$$

$$e_{\text{algebraic}} = \sum_i s_i^2 (\tilde{x}_i - x_i)^2 + s_i^2 (\tilde{y}_i - y_i)^2$$

- s_i is larger for cameras that are further from Q
- Why don't we just divide by s_i ?
 - s_i depends on Q so it is unknown as well.

Linear vs non-linear algorithms

- e_{ideal} can be minimized using non-linear optimization
- Linear algorithms are very fast.
 - Only minimizes an algebraic error.
- We can estimate many things using linear algorithms
 - Triangulation, homography, fundamental matrix, projection matrix
 - They all have the problem that they don't minimize the exact error we desire
- Linear algorithms are acceptable in most cases.
- When high accuracy is desired initialize the non-linear optimization with the linear solution.

Learning objectives

After this lecture you should be able to:

- Derive and explain the epipolar line in computer vision
- Derive and apply the fundamental matrix in computer vision
- Derive and apply the essential matrix in computer vision
- Implement the linear algorithm for triangulation
- Explain the pros and cons of using a linear algorithm

Exercise time!

Camera calibration

Morten R. Hannemose, mohan@dtu.dk

February 23, 2024

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- implement and use the direct linear transformation to calibrate a pinhole camera
- implement and use Zhang's method (2000) to calibrate a pinhole camera using checkerboards

Presentation topics

Direct linear transformation

Zhang's method (2000)

Reprojection error

Non-linear calibration

Practical remarks

Culmination of previous weeks

- Pinhole camera model
- Homogeneous coordinates
- Homographies
- Linear algorithms
- Calibration

Direct linear transformation

Direct linear transformation

Start from the projection equation

$$\mathbf{q}_i = \mathcal{P}Q_i$$

$$\begin{bmatrix} sx_i \\ sy_i \\ s \end{bmatrix} = \mathcal{P} \begin{bmatrix} X_i \\ Y_i \\ Z_i \\ 1 \end{bmatrix}$$

then rearrange into the form $\mathbf{B}^{(i)} \text{flatten}(\mathcal{P}^T) = \mathbf{0}$.

I use the form explained in LN: 2.7 Camera Resection.

Direct linear transformation i

$$\mathbf{q}_i = \mathcal{P}\mathbf{Q}_i,$$

$$\mathbf{0} = \mathbf{q}_i \times \mathcal{P}\mathbf{Q}_i,$$

$$= [\mathbf{q}_i]_{\times} \mathcal{P}\mathbf{Q}_i,$$

$$= \mathbf{B}^{(i)} \text{flatten}(\mathcal{P}^T),$$

where (continue to next slide) ...

Direct linear transformation ii

$$\mathbf{0} = \mathbf{B}^{(i)} \text{flatten}(\mathcal{P}^T),$$

$$\begin{aligned}\mathbf{B}^{(i)} &= \begin{bmatrix} 0 & -X_i & X_i y_i & 0 & -Y_i & Y_i y_i & 0 & -Z_i & Z_i y_i & 0 & -1 & y_i \\ X_i & 0 & -X_i x_i & Y_i & 0 & -Y_i x_i & Z_i & 0 & -Z_i x_i & 1 & 0 & -x_i \\ -X_i y_i & X_i x_i & 0 & -Y_i y_i & Y_i x_i & 0 & -Z_i y_i & Z_i x_i & 0 & -y_i & x_i & 0 \end{bmatrix}, \\ &= \mathbf{Q}_i \otimes [\mathbf{q}_i / s]_{\times}\end{aligned}$$

$$\text{flatten}(\mathcal{P}^T) = [\mathcal{P}_{11} \ \mathcal{P}_{21} \ \mathcal{P}_{31} \ \mathcal{P}_{12} \ \mathcal{P}_{22} \ \mathcal{P}_{32} \ \mathcal{P}_{13} \ \mathcal{P}_{23} \ \mathcal{P}_{33} \ \mathcal{P}_{14} \ \mathcal{P}_{24} \ \mathcal{P}_{34}]^T$$

Direct linear transformation

Now let

$$\mathbf{0} = \mathbf{B} \text{ flatten}(\mathcal{P}^T),$$

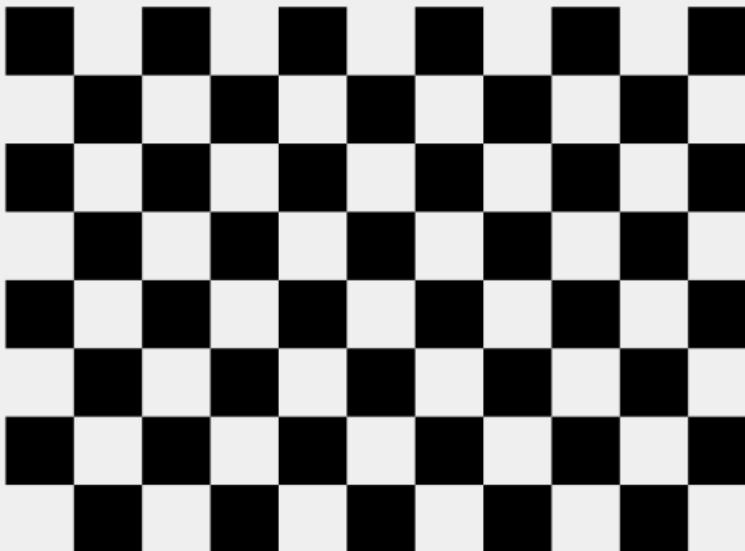
where

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}^{(1)} \\ \mathbf{B}^{(2)} \\ \vdots \end{bmatrix},$$

and solve using SVD on \mathbf{B} .

Zhang's method (2000)

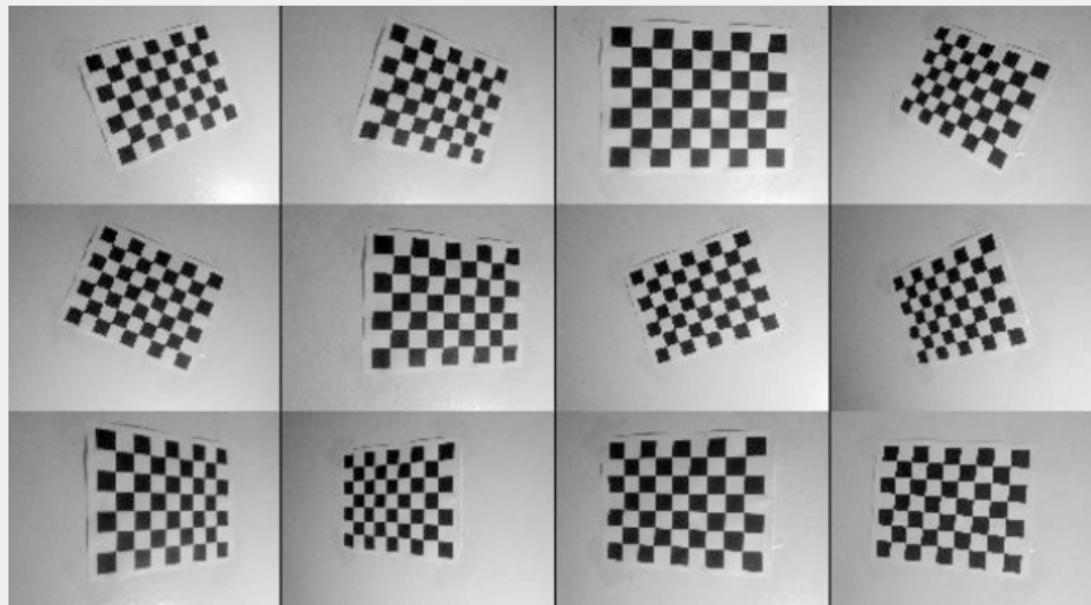
Using a checkerboard



www.calib.io | 8x11 | Checker Size: 15 mm.

Using a checkerboard

View the checkerboard in different poses:



Using a checkerboard

Problem:

Each view i has a different rotation \mathbf{R}_i and translation \mathbf{t}_i

How do we find all \mathcal{P}_i ?

Zhang's method

First, assume all checkerboard corners are in the $Z = 0$ plane:

$$\mathbf{Q}_j = \begin{bmatrix} X_j \\ Y_j \\ 0 \\ 1 \end{bmatrix}$$

Simplifying the projection equation

Let $\mathbf{r}_i^{(c)}$ is the c^{th} column of \mathbf{R}_i . Now projection is

$$\mathbf{q}_{ij} = \mathcal{P}_i \mathbf{Q}_j = \mathbf{K} \begin{bmatrix} \mathbf{r}_i^{(1)} & \mathbf{r}_i^{(2)} & \mathbf{r}_i^{(3)} & \mathbf{t}_i \end{bmatrix} \underbrace{\begin{bmatrix} X_j \\ Y_j \\ 0 \\ 1 \end{bmatrix}}_Q$$

Simplifying the projection equation

Let $\mathbf{r}_i^{(c)}$ is the c^{th} column of \mathbf{R}_i . Now projection is

$$q_{ij} = \mathcal{P}_i Q_j = \mathbf{K} \begin{bmatrix} \mathbf{r}_i^{(1)} & \mathbf{r}_i^{(2)} & \mathbf{r}_i^{(3)} & \mathbf{t}_i \end{bmatrix} \underbrace{\begin{bmatrix} X_j \\ Y_j \\ 0 \\ 1 \end{bmatrix}}_Q$$

$$= \mathbf{K} \begin{bmatrix} \mathbf{r}_i^{(1)} & \mathbf{r}_i^{(2)} & \mathbf{t}_i \end{bmatrix} \underbrace{\begin{bmatrix} X_j \\ Y_j \\ 1 \end{bmatrix}}_{\tilde{Q}_j}.$$

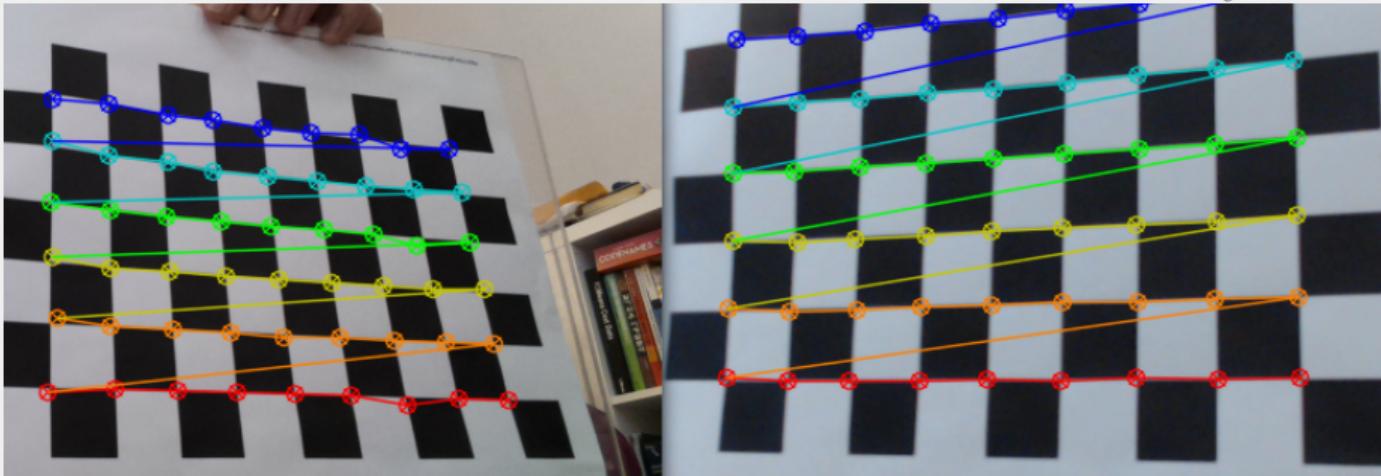
From projections to homographies

$$\begin{aligned} \mathbf{q}_{ij} &= \underbrace{\mathbf{K} \begin{bmatrix} \mathbf{r}_i^{(1)} & \mathbf{r}_i^{(2)} & \mathbf{t}_i \end{bmatrix}}_{\mathbf{H}_i} \underbrace{\begin{bmatrix} X_j \\ Y_j \\ 1 \end{bmatrix}}_{\tilde{\mathbf{Q}}_j} \\ &= \mathbf{H}_i \tilde{\mathbf{Q}}_j \end{aligned}$$

The homographies \mathbf{H}_i can be determined from the plane-plane correspondence $\tilde{\mathbf{Q}}_j$ to \mathbf{q}_{ij} (week 2).

Corner correspondences

Need to find the **unique** positions of corners q_{ij} .



**Find all H_i from corners \tilde{Q}_j and
projections q_{ij} with $q_{ij} = H_i \tilde{Q}_j$**

For example, using SVD.

From homographies to the camera matrix

$$\mathbf{H}_i = \begin{bmatrix} \mathbf{h}_i^{(1)} & \mathbf{h}_i^{(2)} & \mathbf{h}_i^{(3)} \end{bmatrix} = \lambda_i \mathbf{K} \begin{bmatrix} \mathbf{r}_i^{(1)} & \mathbf{r}_i^{(2)} & \mathbf{t}_i \end{bmatrix}.$$

$\mathbf{r}_i^{(1)}$ and $\mathbf{r}_i^{(2)}$ are orthonormal, i.e.

$$\mathbf{r}_i^{(1)\top} \mathbf{r}_i^{(1)} = \mathbf{r}_i^{(2)\top} \mathbf{r}_i^{(2)} = 1,$$

$$\mathbf{r}_i^{(1)\top} \mathbf{r}_i^{(2)} = \mathbf{r}_i^{(2)\top} \mathbf{r}_i^{(1)} = 0.$$

From homographies to the camera matrix

Express $r_i^{(\alpha)}$ using $h_i^{(\alpha)}$:

$$h_i^{(\alpha)} = \lambda_i K r_i^{(\alpha)}, \Leftrightarrow$$

$$K^{-1} h_i^{(\alpha)} = \lambda_i r_i^{(\alpha)}.$$

From homographies to the camera matrix

Express $\mathbf{r}_i^{(\alpha)}$ using $\mathbf{h}_i^{(\alpha)}$:

$$\mathbf{h}_i^{(\alpha)} = \lambda_i \mathbf{K} \mathbf{r}_i^{(\alpha)}, \Leftrightarrow$$

$$\mathbf{K}^{-1} \mathbf{h}_i^{(\alpha)} = \lambda_i \mathbf{r}_i^{(\alpha)}.$$

Now the constraints from the previous slide are:

$$\mathbf{h}_i^{(1)\top} \mathbf{K}^{-\top} \mathbf{K}^{-1} \mathbf{h}_i^{(2)} = 0,$$

$$\mathbf{h}_i^{(1)\top} \mathbf{K}^{-\top} \mathbf{K}^{-1} \mathbf{h}_i^{(1)} = \mathbf{h}_i^{(2)\top} \mathbf{K}^{-\top} \mathbf{K}^{-1} \mathbf{h}_i^{(2)} = \lambda_i^2.$$

We have found constraints on the camera matrix! 😊

Number of constraints

- Two constraints doesn't seem that impressive?

Number of constraints

- Two constraints doesn't seem that impressive?
- Homography has eight degrees of freedom
- Pose of checkerboard has six (3 rotation, 3 translation)
- A single homography can only fix two degrees of freedom of a camera matrix.

Define some new variables i

How to put into practice?

Define the matrix:

$$\mathbf{B} = \mathbf{K}^{-\top} \mathbf{K}^{-1} = \begin{bmatrix} B_{11} & B_{12} & B_{13} \\ B_{12} & B_{22} & B_{23} \\ B_{13} & B_{23} & B_{33} \end{bmatrix},$$
$$\mathbf{b} = [B_{11} \ B_{12} \ B_{22} \ B_{13} \ B_{23} \ B_{33}]^{\top},$$

Define some new variables ii

Now define $\mathbf{v}_i^{(\alpha\beta)}$ such that

$$\mathbf{v}_i^{(\alpha\beta)} \mathbf{b} = \mathbf{h}_i^{(\alpha)T} \mathbf{B} \mathbf{h}_i^{(\beta)}.$$

Then $\mathbf{v}_i^{(\alpha\beta)}$ must be an 1×6 vector given by

$$\begin{aligned}\mathbf{v}_i^{(\alpha\beta)} = & [H_i^{(1\alpha)} H_i^{(1\beta)}, \quad H_i^{(1\alpha)} H_i^{(2\beta)} + H_i^{(2\alpha)} H_i^{(1\beta)}, \quad H_i^{(2\alpha)} H_i^{(2\beta)}, \quad \dots \\ & H_i^{(3\alpha)} H_i^{(1\beta)} + H_i^{(1\alpha)} H_i^{(3\beta)}, \quad H_i^{(3\alpha)} H_i^{(2\beta)} + H_i^{(2\alpha)} H_i^{(3\beta)}, \quad H_i^{(3\alpha)} H_i^{(3\beta)}],\end{aligned}$$

Define some new variables iii

where $H_i^{(rc)}$ is the element in row r and column c of H_i .

From homographies to the camera matrix

Recall our constraints:

$$\begin{aligned} \mathbf{h}_i^{(1)\top} \mathbf{K}^{-\top} \mathbf{K}^{-1} \mathbf{h}_i^{(1)} &= \mathbf{h}_i^{(2)\top} \mathbf{K}^{-T} \mathbf{K}^{-1} \mathbf{h}_i^{(2)} \Leftrightarrow \\ \mathbf{h}_i^{(1)\top} \mathbf{B} \mathbf{h}_i^{(1)} - \mathbf{h}_i^{(2)\top} \mathbf{B} \mathbf{h}_i^{(2)} &= 0 = \\ \left(\mathbf{v}_i^{(11)} - \mathbf{v}_i^{(22)} \right) \mathbf{b} &= 0. \end{aligned}$$

Now we can express our constraints in matrix-form

$$\begin{bmatrix} \mathbf{v}_i^{(12)} \\ \mathbf{v}_i^{(11)} - \mathbf{v}_i^{(22)} \end{bmatrix} \mathbf{b} = \mathbf{0}$$

From homographies to the camera matrix

For all checkerboard poses

$$V\mathbf{b} = \begin{bmatrix} \mathbf{v}_1^{(12)} \\ \mathbf{v}_1^{(11)} - \mathbf{v}_1^{(22)} \\ \mathbf{v}_2^{(12)} \\ \mathbf{v}_2^{(11)} - \mathbf{v}_2^{(22)} \\ \vdots \end{bmatrix} \mathbf{b} = \mathbf{0}$$

When \mathbf{b} is found then \mathbf{K} can be determined using the formulas in Zhang's paper.

Find the camera matrix K through b using $Vb = 0$ and SVD, where V is built from the homographies H_i .

Reprojection error

Is it a good camera calibration? How to find out?

Reprojection error

Is it a good camera calibration? How to find out?

Reproject the points from the checkerboards to the camera, and compare to where we've seen them

$$\tilde{\mathbf{q}}_{ij} = \mathbf{K} [\mathbf{R}_i \ t_i] \mathbf{Q}_j$$

$$\Pi(\tilde{\mathbf{q}}_{ij}) - \Pi(\mathbf{q}_{ij})$$

Reprojection error

We can now compute the reprojection error as the root mean squared error (RMSE)

$$\sqrt{\frac{1}{n} \sum_i \sum_j \|\Pi(\tilde{\mathbf{q}}_{ij}) - \Pi(\mathbf{q}_{ij})\|_2^2}$$

where n is the total number of points.

We have \mathbf{K} , but we are still missing something.

Reprojection error

We can now compute the reprojection error as the root mean squared error (RMSE)

$$\sqrt{\frac{1}{n} \sum_i \sum_j \|\Pi(\tilde{\mathbf{q}}_{ij}) - \Pi(\mathbf{q}_{ij})\|_2^2}$$

where n is the total number of points.

We have \mathbf{K} , but we are still missing something.

How do we recover \mathbf{R}_i and \mathbf{t}_i ?

From homographies to poses

Recall: $\mathbf{H}_i = \begin{bmatrix} \mathbf{h}_i^{(1)} & \mathbf{h}_i^{(2)} & \mathbf{h}_i^{(3)} \end{bmatrix} = \lambda_i \mathbf{K} \begin{bmatrix} \mathbf{r}_i^{(1)} & \mathbf{r}_i^{(2)} & \mathbf{t}_i \end{bmatrix}$

Now we can recover \mathbf{R}_i and \mathbf{t}_i :

$$\mathbf{r}_i^{(1)} = \frac{1}{\lambda_i} \mathbf{K}^{-1} \mathbf{h}_i^{(1)},$$

$$\mathbf{r}_i^{(2)} = \frac{1}{\lambda_i} \mathbf{K}^{-1} \mathbf{h}_i^{(2)},$$

$$\mathbf{t}_i = \frac{1}{\lambda_i} \mathbf{K}^{-1} \mathbf{h}_i^{(3)},$$

$$\mathbf{r}_i^{(3)} = \mathbf{r}_i^{(1)} \times \mathbf{r}_i^{(2)},$$

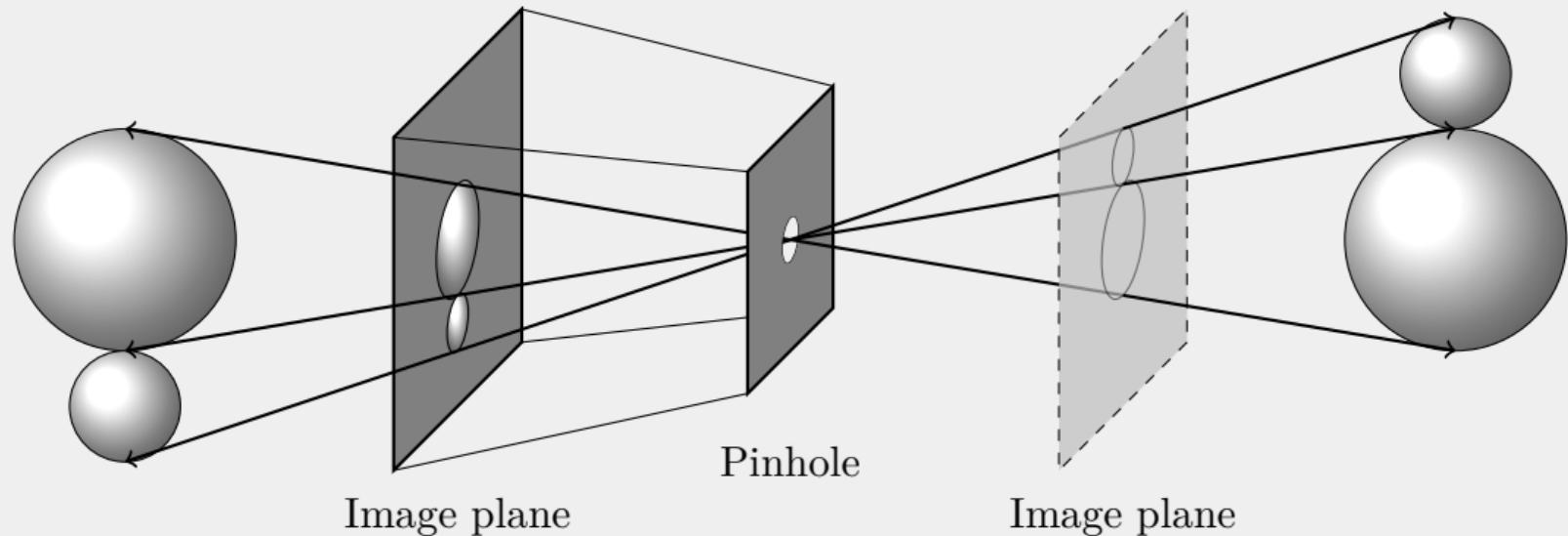
where $\lambda_i = \left\| \mathbf{K}^{-1} \mathbf{h}_i^{(1)} \right\|_2 = \left\| \mathbf{K}^{-1} \mathbf{h}_i^{(2)} \right\|_2$.

Homography and poses

$$\boldsymbol{t}_i = \frac{1}{\lambda_i} \boldsymbol{K}^{-1} \boldsymbol{h}_i^{(3)}$$

What happens if $t_{iz} < 0$? And what does that mean?
Can we ensure correctness?

Pinhole camera revisited



Points behind the camera also get projected to the image plane.
Multiply all 3D points by -1 and they still project to the same place.

Homography and poses

- If $t_{iz} < 0$ the checkerboard is behind the camera.
- How to get the correct solution?
- Estimate \mathbf{R}_i and \mathbf{t}_i again using $-\mathbf{H}_i$ (flipped sign).

Non-linear calibration

Least-squares method

With projection equation

$$\mathbf{q}_{ij} = \mathbf{K} [\mathbf{R}_i \ t_i] \mathbf{Q}_j$$

then

$$E = \sum_{i,j} \left\| \text{dist}\left(\pi\left(\mathbf{K} [\mathbf{R}_i \ t_i] \mathbf{Q}_j\right)\right) - \pi(\mathbf{q}_{ij}) \right\|^2$$

Solution: minimize E w.r.t. \mathbf{K} , \mathbf{R}_i , t_i and lens distortion.

Practical remarks

What images should we take?

- Our two constraints per image are based on $r_i^{(1)}$ and $r_i^{(2)}$.
- Two parallel checkerboards express the same constraints.

What images should we take?

- Our two constraints per image are based on $r_i^{(1)}$ and $r_i^{(2)}$.
- Two parallel checkerboards express the same constraints.
 - at least when we don't consider lens distortion
- Make sure to rotate the checkerboards so they are not parallel.
- Make the checkerboards take up as much of the frame as possible!

How many images?

- Without lens distortion: In theory at least three.
- In practice: **It depends...**
 - Some people use 2000+ images for a single calibration¹
 - Look at the reprojection error
 - Both of the training set and the validation set

Learning objectives

After this lecture you should be able to:

- implement and use the direct linear transformation to calibrate a pinhole camera
- implement and use Zhang's method (2000) to calibrate a pinhole camera using checkerboards

Last week's exercise

The goal of exercise 3.6-3.10 last week was to:

- Understand how to handle when no camera has $R = I$ and $t = 0$.
- Easily find epipolar lines in image 1 from points in image 2.

I have uploaded a new `TwoImageDataCar.npy`, where both have nontrivial \mathbf{R} and \mathbf{t} .

For exercise 3.10, you can use the transpose.

$$\mathbf{q}_1^\top \mathbf{F} \mathbf{q}_2 = 0$$

$$\mathbf{q}_2^\top \mathbf{F}^\top \mathbf{q}_1 = 0$$

If you have already solved the exercise, loading the new file should not take long

Quiz & exercise time!

Nonlinear optimization

camera calibration and triangulation

Morten R. Hannemose, mohan@dtu.dk

March 7, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



Welcome back to myself

About me: Morten Rieger Hannemose
Assistant prof. at DTU Compute,
Section for Visual Computing
Background in computer vision and
differentiable rendering

I work with various applications within
computer vision
camera calibration, 3D scanning, human
pose estimation, and skin cancer.



Misc info

- Weekly quizzes
 - From 2023 exam
 - Do after weekly exercise
 - You can still do previous weeks.
- Let me know if you have any questions
 - About the course
 - About what I'm saying
 - Etc.

**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- perform non-linear optimization within computer vision
- explain the principle behind Levenberg-Marquardt
- compute the Jacobian of a function
- reason about different parameterizations of rotations in optimization

Presentation topics

Non-linear least-squares optimization

- Levenberg–Marquardt

Gradients

- Analytical gradients

- Finite differences approximation

Rotations in optimization

Camera calibration

Non-linear least-squares optimization

Least-squares problems

Computer vision has many problems that can be written as

$$\min_{\boldsymbol{x}} \|g(\boldsymbol{x}) - \boldsymbol{y}\|_2^2.$$

Make all parameters into a vector \boldsymbol{x} and optimize everything!

Least-squares problems

Computer vision has many problems that can be written as

$$\min_{\boldsymbol{x}} \|g(\boldsymbol{x}) - \boldsymbol{y}\|_2^2.$$

Make all parameters into a vector \boldsymbol{x} and optimize everything!

- Homography estimation
- Pose estimation
- Bundle adjustment
- Camera calibration with and without lens distortion
- Triangulation
- ...

Reformulate a bit

$$\begin{aligned} e(\mathbf{x}) &= \left\| \underbrace{g(\mathbf{x}) - \mathbf{y}}_{f(\mathbf{x})} \right\|_2^2 \\ &= \|f(\mathbf{x})\|_2^2 \\ &= f(\mathbf{x})^\top f(\mathbf{x}) \end{aligned}$$

Levenberg-Marquardt

We solve the problem in an iterative fashion. At the k^{th} iteration:

- Our current guess of \boldsymbol{x} is \boldsymbol{x}_k
- Replace with f with first order approximation around \boldsymbol{x}_k

$$f(\boldsymbol{x}_k + \boldsymbol{\delta}) \approx f(\boldsymbol{x}_k) + \mathbf{J}\boldsymbol{\delta}$$

- \mathbf{J} is the Jacobian that contains all first order derivatives of f at \boldsymbol{x}_k .

Levenberg-Marquardt

We solve the problem in an iterative fashion. At the k^{th} iteration:

- Our current guess of \boldsymbol{x} is \boldsymbol{x}_k
- Replace with f with first order approximation around \boldsymbol{x}_k

$$f(\boldsymbol{x}_k + \boldsymbol{\delta}) \approx f(\boldsymbol{x}_k) + \mathbf{J}\boldsymbol{\delta}$$

- \mathbf{J} is the Jacobian that contains all first order derivatives of f at \boldsymbol{x}_k .

How are the the Jacobians of f and g related?

Levenberg-Marquardt

We solve the problem in an iterative fashion. At the k^{th} iteration:

- Our current guess of \boldsymbol{x} is \boldsymbol{x}_k
- Replace with f with first order approximation around \boldsymbol{x}_k

$$f(\boldsymbol{x}_k + \boldsymbol{\delta}) \approx f(\boldsymbol{x}_k) + \mathbf{J}\boldsymbol{\delta}$$

- \mathbf{J} is the Jacobian that contains all first order derivatives of f at \boldsymbol{x}_k .

How are the the Jacobians of f and g related?

They are **identical!**

The Jacobian

f is defined as

$$f(\mathbf{x}) = f \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix}$$

The Jacobian of f is given by

$$\mathbf{J} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Levenberg-Marquardt

The sum of squared errors is then:

$$e(\mathbf{x}_k + \boldsymbol{\delta}) = \|f(\mathbf{x}_k + \boldsymbol{\delta})\|_2^2 \approx$$

$$\tilde{e}(\mathbf{x}_k + \boldsymbol{\delta}) = (f(\mathbf{x}_k) + \mathbf{J}\boldsymbol{\delta})^\top (f(\mathbf{x}_k) + \mathbf{J}\boldsymbol{\delta})$$

Levenberg-Marquardt

The sum of squared errors is then:

$$e(\mathbf{x}_k + \boldsymbol{\delta}) = \|f(\mathbf{x}_k + \boldsymbol{\delta})\|_2^2 \approx$$

$$\begin{aligned}\tilde{e}(\mathbf{x}_k + \boldsymbol{\delta}) &= (f(\mathbf{x}_k) + \mathbf{J}\boldsymbol{\delta})^\top (f(\mathbf{x}_k) + \mathbf{J}\boldsymbol{\delta}) \\ &= f(\mathbf{x}_k)^\top f(\mathbf{x}_k) + 2(\mathbf{J}\boldsymbol{\delta})^\top f(\mathbf{x}_k) + (\mathbf{J}\boldsymbol{\delta})^\top \mathbf{J}\boldsymbol{\delta}^\top \\ &= f(\mathbf{x}_k)^\top f(\mathbf{x}_k) + 2\boldsymbol{\delta}^\top \mathbf{J}^\top f(\mathbf{x}_k) + \boldsymbol{\delta}^\top \mathbf{J}^\top \mathbf{J}\boldsymbol{\delta}^\top\end{aligned}$$

which is a second order approximation of e using only first order derivatives of f 😊

Levenberg-Marquardt

We can minimize our approximation \tilde{e} instead of e .

Take the derivative of $\tilde{e}(\mathbf{x}_k + \boldsymbol{\delta})$

$$\begin{aligned}\tilde{e}(\mathbf{x}_k + \boldsymbol{\delta}) &= f(\mathbf{x}_k)^\top f(\mathbf{x}_k) + 2\boldsymbol{\delta}^\top \mathbf{J}^\top f(\mathbf{x}_k) + \boldsymbol{\delta}^\top \mathbf{J}^\top \mathbf{J} \boldsymbol{\delta}^\top \\ \frac{\partial \tilde{e}(\mathbf{x}_k + \boldsymbol{\delta})}{\partial \boldsymbol{\delta}} &= 2\mathbf{J}^\top f(\mathbf{x}_k) + 2\mathbf{J}^\top \mathbf{J} \boldsymbol{\delta}.\end{aligned}$$

Levenberg-Marquardt

Find the optimum by setting the derivative equal to zero

$$\begin{aligned} 2\mathbf{J}^T f(\mathbf{x}_k) + 2\mathbf{J}^T \mathbf{J}\boldsymbol{\delta} = \mathbf{0} &\Leftrightarrow \\ \mathbf{J}^T \mathbf{J}\boldsymbol{\delta} &= -\mathbf{J}^T f(\mathbf{x}_k). \end{aligned}$$

Solve for $\boldsymbol{\delta}$, and then set $\mathbf{x}_{k+1} = \mathbf{x}_k + \boldsymbol{\delta}$.

Rinse and repeat!

The λ parameter

The approximation \tilde{e} is better the closer we are to the minimum.

When far away, it can be less good. What can we do in this case?

The λ parameter

The approximation \tilde{e} is better the closer we are to the minimum.

When far away, it can be less good. What can we do in this case?

Fall back to **gradient descent**!

$$(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I})\boldsymbol{\delta} = -\mathbf{J}^\top f(\mathbf{x}_k).$$

Decrease λ when $e(\mathbf{x}_{k+1}) < e(\mathbf{x}_k)$, otherwise increase it.

Levenberg-Marquardt summary

Levenberg-Marquardt is often well suited to the types of problems we encounter in computer vision.

It exploits the nature of the least-squares problem to get a **second order** method, using only first order derivatives.

Gradients

Gradients

How do we compute J ?

Gradients

How do we compute \mathbf{J} ?

Recall that \mathbf{J} is just a lot of derivatives in a matrix.

- Analytical gradients
- Automatic differentiation
 - Reverse mode automatic differentiation (backpropagation)
 - Forward mode automatic differentiation (dual numbers)
- Finite differences approximation

Gradients

How do we compute J ?

Recall that J is just a lot of derivatives in a matrix.

- Analytical gradients
- Automatic differentiation
 - Reverse mode automatic differentiation (backpropagation)
 - Forward mode automatic differentiation (dual numbers)
- Finite differences approximation

Analytical gradients and where to find them?

Analytical gradients and where to find them?

Get out pen and paper (or Maple or Sympy) and find the exact derivative.

Analytical gradients - triangulation example

The reprojection error of a point in 3D in n cameras

$$f(\mathbf{Q}) = \begin{bmatrix} \Pi(\mathcal{P}_1 \Pi^{-1}(\mathbf{Q})) - \tilde{\mathbf{q}}_1 \\ \vdots \\ \Pi(\mathcal{P}_n \Pi^{-1}(\mathbf{Q})) - \tilde{\mathbf{q}}_n \end{bmatrix}$$

- f returns a vector of length $2n$
- \mathbf{Q} has three parameters
- thus \mathbf{J} is $2n \times 3$.

Analytical gradients - triangulation example

The projected point in homogeneous coordinates:

$$\begin{aligned} q &= \underbrace{\mathcal{P} \Pi^{-1}(Q)}_{Q_h} \\ &= \begin{bmatrix} sx \\ sy \\ s \end{bmatrix} = \begin{bmatrix} \mathcal{P}^{(1)} \\ \mathcal{P}^{(2)} \\ \mathcal{P}^{(3)} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \end{aligned}$$

Analytical gradients - triangulation example

Let's focus on x

$$x = \frac{\boldsymbol{p}^{(1)} \mathbf{Q}_h}{\boldsymbol{p}^{(3)} \mathbf{Q}_h} = \frac{\mathcal{P}_{11}X + \mathcal{P}_{12}Y + \mathcal{P}_{13}Z + \mathcal{P}_{14}}{\mathcal{P}_{31}X + \mathcal{P}_{32}Y + \mathcal{P}_{33}Z + \mathcal{P}_{34}}$$

A single element of \mathbf{J} is then given by:

$$\frac{d}{dX}x = \frac{\mathcal{P}_{11}}{\boldsymbol{p}^{(3)} \mathbf{Q}_h} - \frac{\mathcal{P}_{31}(\boldsymbol{p}^{(1)} \mathbf{Q}_h)}{(\boldsymbol{p}^{(3)} \mathbf{Q}_h)^2}$$

Analytical gradients - summary

Analytical gradients are extremely useful

- + very fast to compute
- + accurate
- complicated to derive and implement

Many functions in OpenCV return the Jacobian in addition the values themselves.

Finite differences approximation - forward differences

A first order Taylor expansion of f

$$f(x + h) = f(x) + \frac{d}{dx}f(x)h + O(h)$$

Finite differences approximation - forward differences

A first order Taylor expansion of f

$$f(x + h) = f(x) + \frac{d}{dx}f(x)h + O(h) \Leftrightarrow$$

can be rewritten to

$$\frac{d}{dx}f(x) = \frac{f(x + h) - f(x)}{h} + O(h).$$

This is called forward differences. (2-point)

Finite differences approximation - central differences

A second order Taylor expansion of f evaluated at $x - h$ and $x + h$ can be rearranged to yield

$$\frac{d}{dx}f(x) = \frac{f(x+h) - f(x-h)}{2h} + O(h^2).$$

This is called central differences, which is more accurate ($O(h^2)$), but requires two new evaluations of f .

Finite differences approximation - in practice

Principle can only be applied to one parameter at a time.

To compute \mathbf{J} you need to evaluate f once (or twice) for each element in \mathbf{x} (which can be a lot)

Number of evaluations can be reduced if you know the sparsity structure of \mathbf{J} , i.e. which elements of \mathbf{x} affect which elements of $f(\mathbf{x})$

Finite differences approximation - h

How do we choose h ?

A good choice is a fixed percentage of x .

Can this have any issues?

Finite differences approximation - h

How do we choose h ?

A good choice is a fixed percentage of x .

Can this have any issues?

When x is zero, a percentage is also zero, this needs to be handled.

Finite differences approximation - summary

They only give an approximation of the gradients

- + convenient
- lots of (unnecessary) computation
- has parameter h
- numeric problems

Only use when speed and robustness are not your primary concerns.

Short Break

Rotations in optimization

Rotations in optimization

- Rotations are usually 3×3 matrices, but have just three degrees of freedom.
- How can we parametrize rotations?

Rotations in optimization

- Rotations are usually 3×3 matrices, but have just three degrees of freedom.
- How can we parametrize rotations?
 - optimizing all 9 numbers can yield something that is no longer a rotation matrix

Rotations in optimization - Euler angles

Euler angles ($\theta_x, \theta_y, \theta_z$) use only three numbers

- Suffers from gimbal lock, i.e. that two parameters do the same in certain configurations
 - unsuitable for optimization

Rotations in optimization - Euler angles

Euler angles ($\theta_x, \theta_y, \theta_z$) use only three numbers

- Suffers from gimbal lock, i.e. that two parameters do the same in certain configurations
 - unsuitable for optimization
- If the rotation is close to identity, we are far from gimbal lock
 - no problems in this case
 - optimize over a rotation relative to the initial guess
 - $\theta_x, \theta_y, \theta_z$ will usually stay close to zero
 - only if initial guess is good
 - for finite difference start from $\theta_x = \theta_y = \theta_z = 2\pi$

Rotation Explorer

Rotations in optimization - Axis-angle

Axis-angle represents a rotation as a rotation of θ around an axis v , where $\|v\|_2 = 1$.

Store as a vector v/θ (only three elements).

How OpenCV returns rotations (rvec). Can be converted to a rotation matrix with cv2.Rodrigues.

Has singularity at 0, but works well in most cases.

Rotations in optimization - Quaternions

Quaternions are the gold standard for representing rotations.

A quaternion uses four numbers represent a rotation (q_1, q_2, q_3, q_4) subject to $q_1^2 + q_2^2 + q_3^2 + q_4^2 = 1$.

Smooth and suffers from no problems, however we have to normalize the quaternion each optimization step, to ensure it is still a valid quaternion.

Software packages

There are **many** implementations out there!

- Ceres (C++)
 - Dual numbers (no implementing derivatives, 😊)
 - Quaternion parameterization
- `scipy.optimize.least_squares`
 - Easy to use
 - Finite differences or analytical gradients
- Jax

Camera calibration

Outline

More tips and tricks

- Checkerboard alternatives
- Subpixel estimation
- Overfitting
- Bundle adjustment

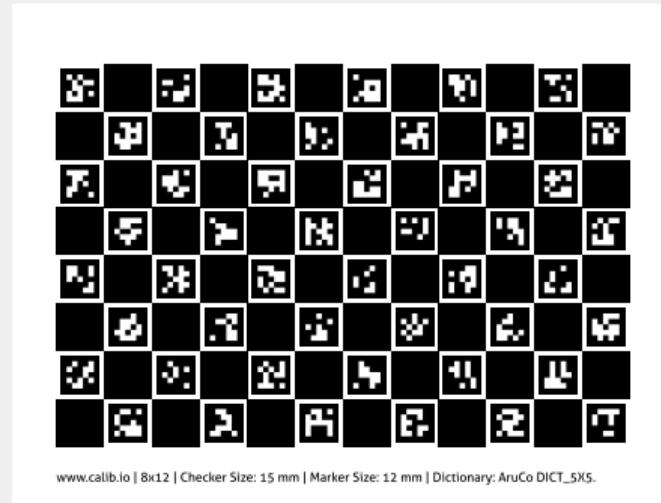
Checkerboard

OpenCV needs to see the all corners of the checkerboard in order to detect it.

Getting the entire image plane covered in detected points is hard, especially near the edges.

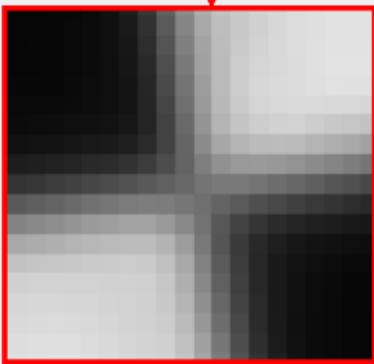
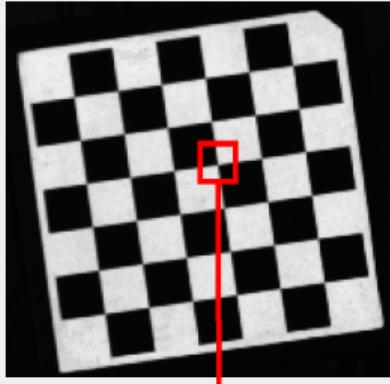
ChArUco

A ChArUco boards is a checkerboard with ArUco markers.



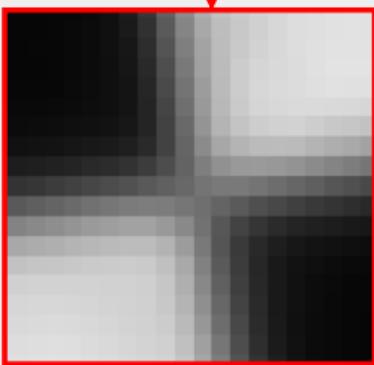
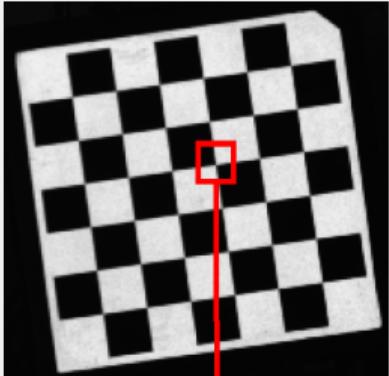
Possible to detect partial checkerboards!

Subpixel corner estimation



- Where is the corner?

Subpixel corner estimation



- Where is the corner?
- Look at a neighbourhood around a corner
- OpenCV has `cv2.cornerSubPix` but it is 💩
- What to do?

Good method for subpixel corner estimation

A good subpixel corner estimator is presented in Schops et al.¹

Generate n random 2D vectors relative to the corner \mathbf{s}_i

$$C_{sym}(\mathbf{H}) = \sum_{i=1}^n \left((I(\mathbf{H}(\mathbf{s}_i)) - I(\mathbf{H}(-\mathbf{s}_i)))^2 \right)$$

\mathbf{H} maps a point using the homography, and I interpolates the value of the image at the given point. Minimize C_{sym} .

Camera calibration - A cautionary tale

- Including more distortion parameters in your model will always give a lower re-projection error

Camera calibration - A cautionary tale

- Including more distortion parameters in your model will always give a lower re-projection error
- ...for the **images used to calibration!**
 - Will not necessarily generalize.
- This can be overfitting
- Use **cross-validation**

Camera calibration - How to do cross validation

- For each checkerboard, split your detected corners into a validation and training set
- Do the calibration using all training corners
- Use the estimated checkerboard pose to reproject the validation corners

Bundle adjustment

- The checkerboard can have imperfections, such as not being flat
 - The checkerboards you get today are not very flat
- Optimize everything again, including the 3D positions of the corners

Learning objectives

After this lecture you should be able to:

- perform non-linear optimization within computer vision
- explain the principle behind Levenberg-Marquardt
- compute the Jacobian of a function
- reason about different parameterizations of rotations in optimization

Info about exercise and exercise time!

Simple features

Morten R. Hannemose, mohan@dtu.dk

March 14, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- explain the image correspondence problem
- explain the use of image features (key points, interest points)
- understand and implement image filtering
- implement Harris corner detection
- run Canny edge detection

Presentation topics

Image correspondence problem

Features and feature descriptors

Simple features

Filtering and convolution

Harris corners

Canny edges

Image correspondence problem

Image correspondence

The problem of matching two parts of different images:

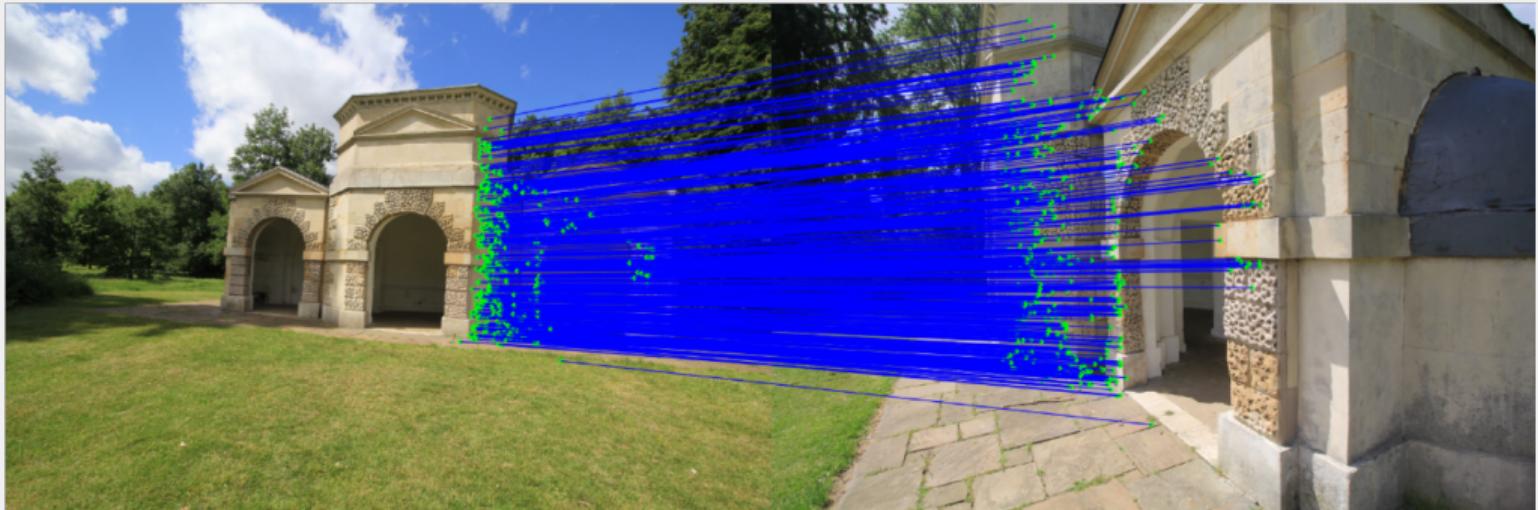


Image correspondence

More images of the same scene

- Correspondence exists between parts that are visible in each image
- Not all corresponding points have a unique pattern
- Idea: Only choose points that can be identified uniquely

What are good points?



What are good points?

We are able to identify the same corners in multiple images



Problems

Movement of camera

- Scale – distance of camera to object
- Rotation – objects and camera is rotated between frames
- Translation – movement of camera from one place to another

Result: Appearance change dependent on distance to camera

Problems

Movement of camera



Problems

Other issues

- Occlusion – objects can be in front of other objects
- Lighting change – darker/lighter, color of light (change in spectrum), direction of light
- Clutter – many objects in a scene

Solution

Invariance to imaging problems

- Focus on points and a small area around each point
- Two aspects
 - identify key points (focus of today)
 - characterize pattern around point (week 8)

Features and feature descriptors

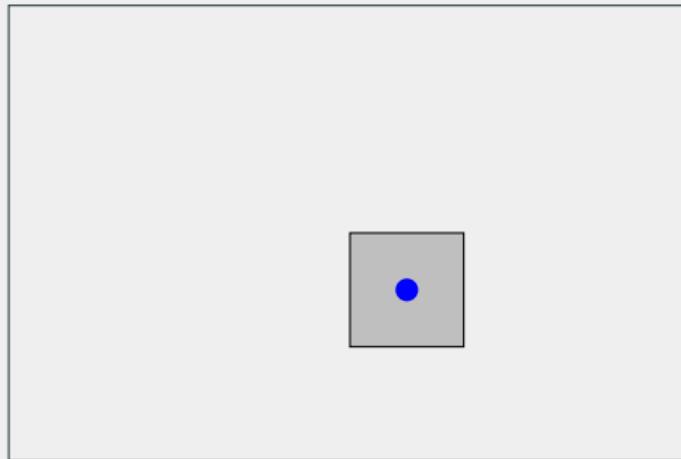
First some terminology:

Key points, interest points, and feature points are the same – namely points in an image with a coordinate position (also between pixels).

Descriptors, key point descriptors, interest point descriptors, and feature descriptors are the same and used for characterizing the pattern around a point. Usually a vector that can be matched between images.

The image patch feature descriptor

One easy way of describing a feature is to use the local pixel values around the feature.



The image patch matching

Stretch the image patch into a vector \mathbf{f} , and that is your descriptor.

A simple comparison operator $d(\mathbf{f}_1, \mathbf{f}_2)$ just uses the scalar product as the distance between features:

$$d(\mathbf{f}_1, \mathbf{f}_2) = \mathbf{f}_1^\top \mathbf{f}_2$$

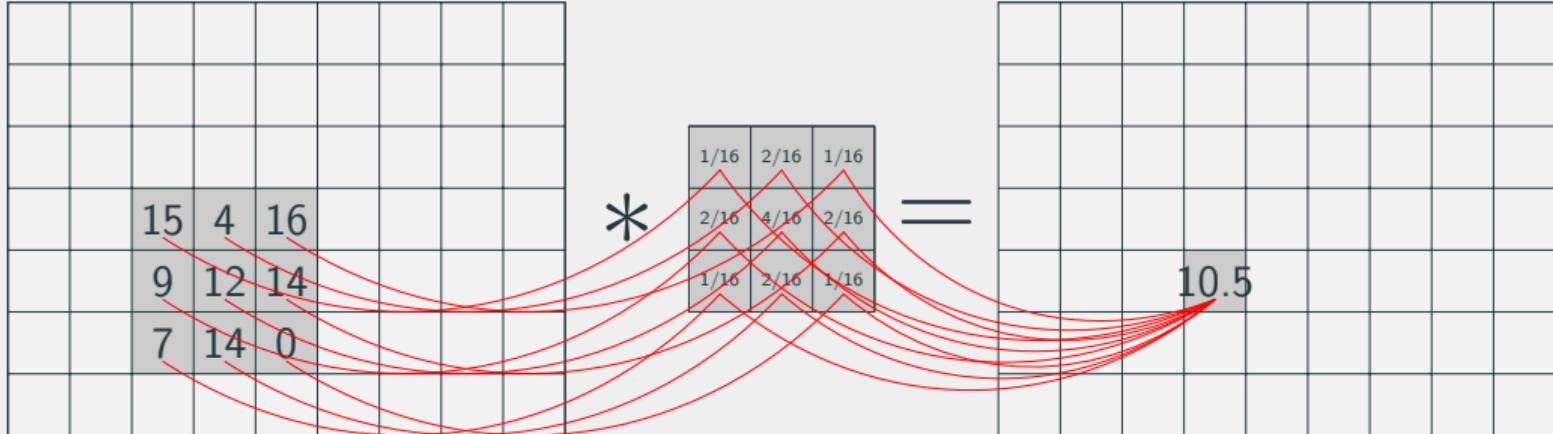
Good with only small transformations, low noise, and unchanged environment.

More robust (but more complicated) feature descriptors will follow in two weeks.

Simple features

Filtering and convolution

$$I(x; t) = \int_{\xi \in \mathbb{R}^2} I(x - \xi) g(\xi; t) d\xi \approx \sum_{i=-k}^k \sum_{j=-k}^k I(x - i, y - j) g(i, j)$$



2D convolution

We use filtering and convolution for the same operation and use the symbol $*$ for convolution. Convolution is commutative

$$I_g = g * I = I * g$$

2D Gaussian $g : \mathbb{R}^2 \times \mathbb{R}_+ \rightarrow \mathbb{R}$

$$g(x, y; \sigma^2) = \frac{1}{2\pi\sigma^2} \exp\left(\frac{-(x^2 + y^2)}{2\sigma^2}\right). \quad (1)$$

Convolution – separability

Useful property: The (isotropic) Gaussian is separable

$$g(x, y; \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-x^2}{2\sigma^2}\right) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-y^2}{2\sigma^2}\right)$$

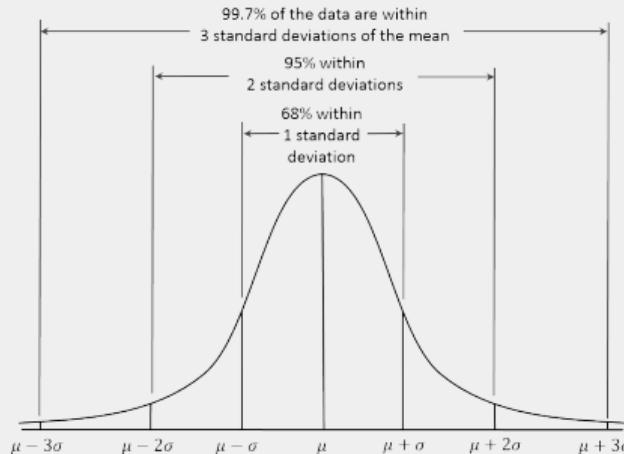
This separability means that

$$I_g = (\mathbf{g} * \mathbf{g}^T) * I = \mathbf{g} * (\mathbf{g}^T * I)$$

where \mathbf{g} is a vector of the Gaussian.

Size of Gaussian filter

- Ideally: Infinitely big – but this is impractical
- Empiric rule: 3σ , 4σ or 5σ rule – example:
 - if $\sigma = 2$, filter size of Gaussian is $2 \cdot 5 \cdot 2 + 1 = 21$
 - if $\sigma = 20$, filter size of Gaussian is $2 \cdot 5 \cdot 20 + 1 = 201$



Demo – Gaussian filtering

cv2.sepFilter2D

Derivatives of an image

How can we find the derivative of an image in e.g. the x -direction?
What does it mean to take a derivative?

Derivatives of an image

How can we find the derivative of an image in e.g. the x -direction?

What does it mean to take a derivative?

An image is only discrete values.

If we can make the image **continuous** all our math works again

Derivatives of an image using Gaussians

Let \mathbf{g} be column vector of a Gaussian. Consider the blurred image

$$I_b = \mathbf{g} * \mathbf{g}^T I.$$

Derivatives of an image using Gaussians

Let \mathbf{g} be column vector of a Gaussian. Consider the blurred image

$$I_b = \mathbf{g} * \mathbf{g}^\top I.$$

Then the derivative of I_b in the x -direction is

$$\begin{aligned}\frac{\partial}{\partial x} I_b &= \frac{\partial}{\partial x} (\mathbf{g} * \mathbf{g}^\top * I) \\ &= \mathbf{g} * \left(\frac{\partial}{\partial x} \mathbf{g}^\top \right) * I \\ &= \mathbf{g} * \mathbf{g}_d^\top * I,\end{aligned}$$

Derivative of the Gaussian

Recall the one-dimensional Gaussian

$$g(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-x^2}{2\sigma^2}\right).$$

We can then compute the derivative

$$g_d(x) = \frac{d}{dx}g(x) = \frac{-x}{\sigma^2}g(x).$$

This is straightforward to implement on a computer.

Short break

Harris corners

Corners (and blobs) are great features because they are easy to describe and detect.

Let's see how much the intensity changes, for a small shift Δ_x, Δ_y

$$\Delta I(x, y, \Delta_x, \Delta_y) = I(x, y) - I(x + \Delta_x, y + \Delta_y)$$

Harris corners are detected as points with *locally maximum* change from a small shift.

A corner is then a local area where $\Delta I(x, y, \Delta_x, \Delta_y)^2$ is large no matter how we choose Δ_x, Δ_y .

Harris corner measure

To improve robustness towards noise the “cornerness” should be true for a local area and not just a point. We apply a Gaussian blur to the above measure.

The measure to check for corners is then

$$\begin{aligned} c(x, y, \Delta_x, \Delta_y) &= g * \Delta I(x, y, \Delta_x, \Delta_y)^2 \\ &= g * (I(x, y) - I(x + \Delta_x, y + \Delta_y))^2, \end{aligned}$$

where $g * \dots$ is the convolution with the Gaussian.

Replace with Taylor approximation

$$I(x + \Delta_x, y + \Delta_y) \approx I(x, y) - \frac{\partial I}{\partial x}(x, y) \Delta_x - \frac{\partial I}{\partial y}(x, y) \Delta_y$$

Replace with Taylor approximation

$$I(x + \Delta_x, y + \Delta_y) \approx I(x, y) - \underbrace{\frac{\partial I}{\partial x}(x, y) \Delta_x}_{I_x} - \underbrace{\frac{\partial I}{\partial y}(x, y) \Delta_y}_{I_y}$$

Replace with Taylor approximation

$$\begin{aligned} I(x + \Delta_x, y + \Delta_y) &\approx I(x, y) - \underbrace{\frac{\partial I}{\partial x}(x, y) \Delta_x}_{I_x} - \underbrace{\frac{\partial I}{\partial y}(x, y) \Delta_y}_{I_y} \\ &= I(x, y) - [I_x \quad I_y] \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix} \end{aligned}$$

I_x and I_y are also values at (x, y) , but we omit it for readability.

Harris corner measure derivation

The corner checking measure is then

$$\begin{aligned} c(x, y, \Delta_x, \Delta_y) &= g * \left(I(x, y) - I(x + \Delta_x, y + \Delta_y) \right)^2 \\ &\approx g * \left(\begin{bmatrix} I_x & I_y \end{bmatrix} \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix} \right)^2 \end{aligned}$$

Harris corner measure derivation

The corner checking measure is then

$$\begin{aligned} c(x, y, \Delta_x, \Delta_y) &= g * \left(I(x, y) - I(x + \Delta_x, y + \Delta_y) \right)^2 \\ &\approx g * \left(\begin{bmatrix} I_x & I_y \end{bmatrix} \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix} \right)^2 \\ &= g * \left(\begin{bmatrix} \Delta_x & \Delta_y \end{bmatrix} \begin{bmatrix} I_x \\ I_y \end{bmatrix} \begin{bmatrix} I_x & I_y \end{bmatrix} \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix} \right) \\ &= g * \left(\begin{bmatrix} \Delta_x & \Delta_y \end{bmatrix} \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix} \right) \end{aligned}$$

Harris corner measure derivation

Finally, we take the convolution with the Gaussian inside

$$\begin{aligned} c(x, y, \Delta_x, \Delta_y) &\approx g * \left([\Delta_x \quad \Delta_y] \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix} \right) \\ &= [\Delta_x \quad \Delta_y] \underbrace{\begin{bmatrix} g * (I_x^2) & g * (I_x I_y) \\ g * (I_x I_y) & g * (I_y^2) \end{bmatrix}}_{C(x,y)} \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix}. \end{aligned}$$

We end up with $C(x, y)$, also known as the **structure tensor**.

The structure tensor

Corners have a large $c(x, y, \Delta_x, \Delta_y)$ regardless of $[\Delta_x \ \Delta_y]$.

$$c(x, y, \Delta_x, \Delta_y) \approx [\Delta_x \ \Delta_y] \mathbf{C}(x, y) \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix}$$

How can we check if this the case from the $\mathbf{C}(x, y)$ matrix?

The structure tensor

Corners have a large $c(x, y, \Delta_x, \Delta_y)$ regardless of $[\Delta_x \ \Delta_y]$.

$$c(x, y, \Delta_x, \Delta_y) \approx [\Delta_x \ \Delta_y] \mathbf{C}(x, y) \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix}$$

How can we check if this the case from the $\mathbf{C}(x, y)$ matrix?

When both eigenvalues of $\mathbf{C}(x, y)$ are large.

The Harris corner metric continued

Let λ_1 and λ_2 be the eigenvalues of $\mathbf{C}(x, y) = \begin{bmatrix} a & c \\ c & b \end{bmatrix}$.

We then have the following

$$\det(\mathbf{C}(x, y)) = \lambda_1\lambda_2 = ab - c^2,$$

$$\text{trace}(\mathbf{C}(x, y)) = \lambda_1 + \lambda_2 = a + b, \text{ let us then define}$$

$$r(x, y) = \det(\mathbf{C}(x, y)) - k \text{ trace}(\mathbf{C}(x, y))^2$$

The Harris corner metric continued

Let λ_1 and λ_2 be the eigenvalues of $\mathbf{C}(x, y) = \begin{bmatrix} a & c \\ c & b \end{bmatrix}$.

We then have the following

$$\det(\mathbf{C}(x, y)) = \lambda_1\lambda_2 = ab - c^2,$$

$$\text{trace}(\mathbf{C}(x, y)) = \lambda_1 + \lambda_2 = a + b, \text{ let us then define}$$

$$\begin{aligned} r(x, y) &= \det(\mathbf{C}(x, y)) - k \text{ trace}(\mathbf{C}(x, y))^2 \\ &= \lambda_1\lambda_2 - k(\lambda_1 + \lambda_2)^2 \\ &= ab - c^2 - k(a + b)^2 \end{aligned}$$

The Harris corner metric

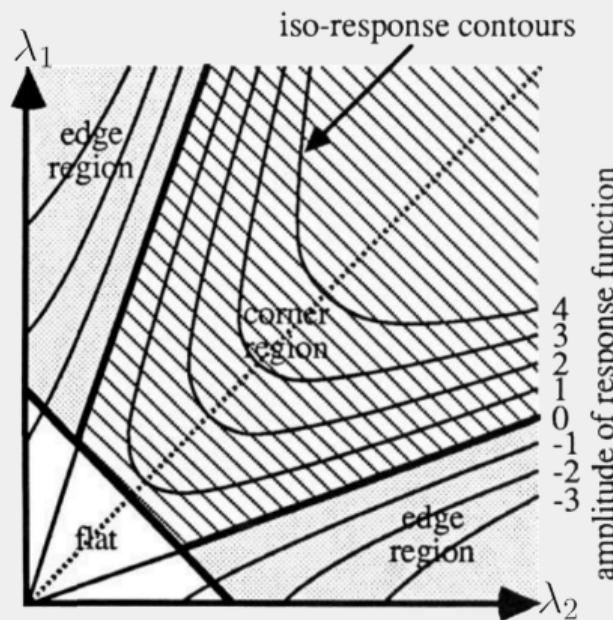
Now we have the Harris corner metric

$$r(x, y) = \lambda_1\lambda_2 - k(\lambda_1 + \lambda_2)^2.$$

k is a free parameter, typically $k = 0.06$.

Notice that $r(x, y)$ is:

- negative for one eigenvalue much greater than the other
- large and positive for large eigenvalues
- small and positive for small eigenvalues



The Harris corner detector

We can now detect corners by finding points where $r(x, y)$ is greater than some threshold τ .

Typically, you can choose τ to be

$$0.1 \cdot \max(r(x, y)) < \tau < 0.8 \cdot \max(r(x, y)).$$

Non-maximum suppression

Find the local maximum of one pixel compared to neighbours

$$(I(x, y) - I(x', y')) > 0 \quad \forall x' \in n(x, y)$$

where $n(x, y)$ is a neighbourhood around the point (x, y)

Suggested procedure

- Initialize a new array **M** with $I(x, y) > \tau$
- Compare original image with one neighbour (e.g to the right)
- Set pixels in **M** to 0 where original not larger than neighbour
 - Repeat for all neighbours
- Find coordinates of all pixels that are 1

Canny edges

Often, we might also consider detecting lines.

To detect (non-straight) lines, we can use the Canny edge detector.

Canny edges

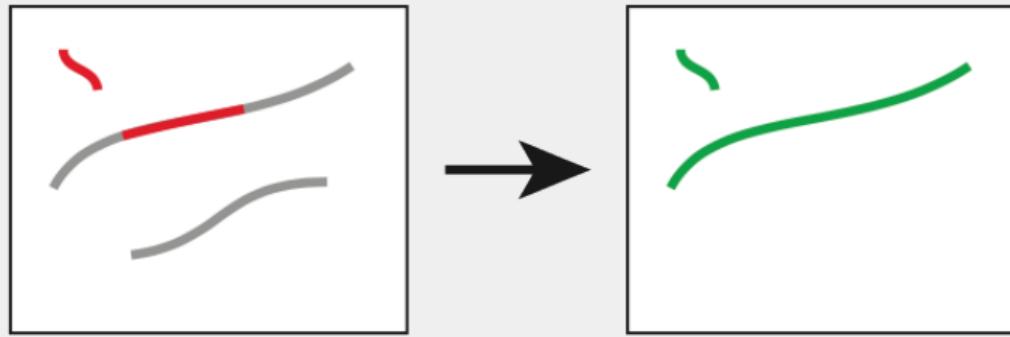
The metric in the canny edge detector is simply the gradient magnitude

$$m(x, y) = \sqrt{I_x^2(x, y) + I_y^2(x, y)}.$$

And the edges are detected by thresholding the magnitude $m(x, y)$, however, in two stages: **seed and grow**.

Canny edges

- Seed: label edges by a threshold where $m(x, y) > \tau_1$
- Grow: with a second threshold where $m(x, y) > \tau_2$, label iff the new points are next to previously labelled edges.
- The threshold values are chosen such that $\tau_1 > \tau_2$



Exercise

During the exercise, you will

- **implement** the Harris corner detector
- **try** the Canny edge detector

Suggestion: After finishing your implementation – try on your own images and perhaps draw some features and see how it works

Learning objectives

After this lecture you should be able to:

- explain the image correspondence problem
- explain the use of image features (key points, interest points)
- understand and implement image filtering
- implement Harris corner detection
- run Canny edge detection

Exercise time!

Robust Model Fitting

Morten R. Hannemose, mohan@dtu.dk

March 21, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- explain how the Hough transform works
- understand and implement RANSAC

Presentation topics

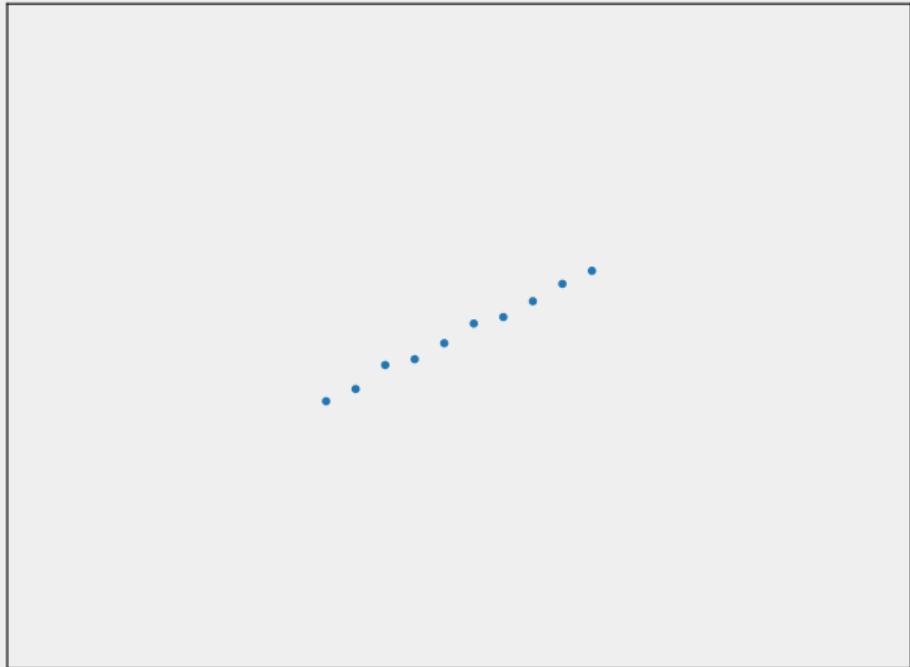
Hough Transform

RANSAC

Recap of lines in homogeneous coordinates

Fitting models

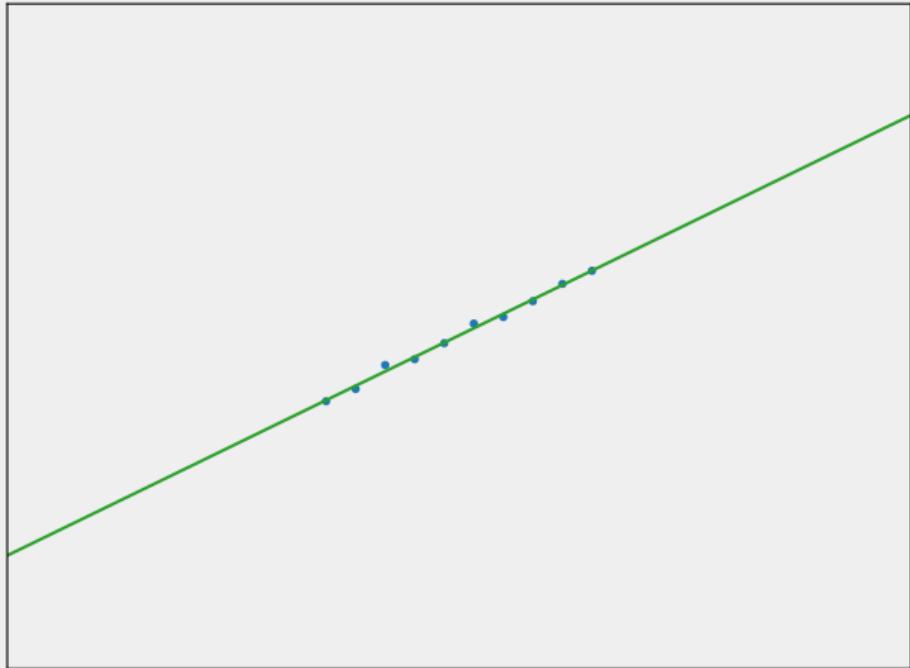
Can we fit a straight line?



Fitting models

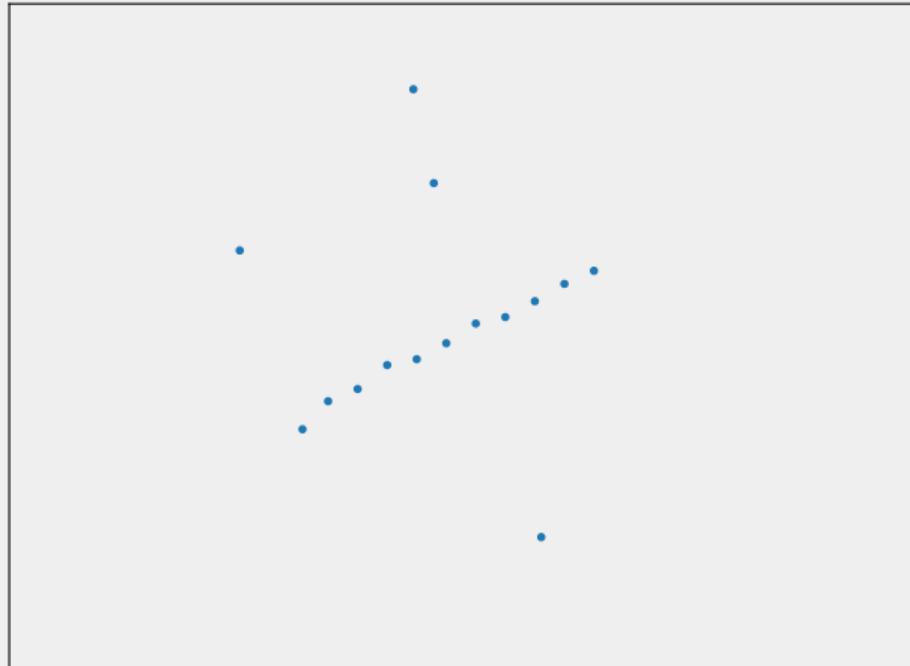
Can we fit a straight line?

Yes we can!



Fitting models

Can we fit a straight line
when there are a few
outliers?

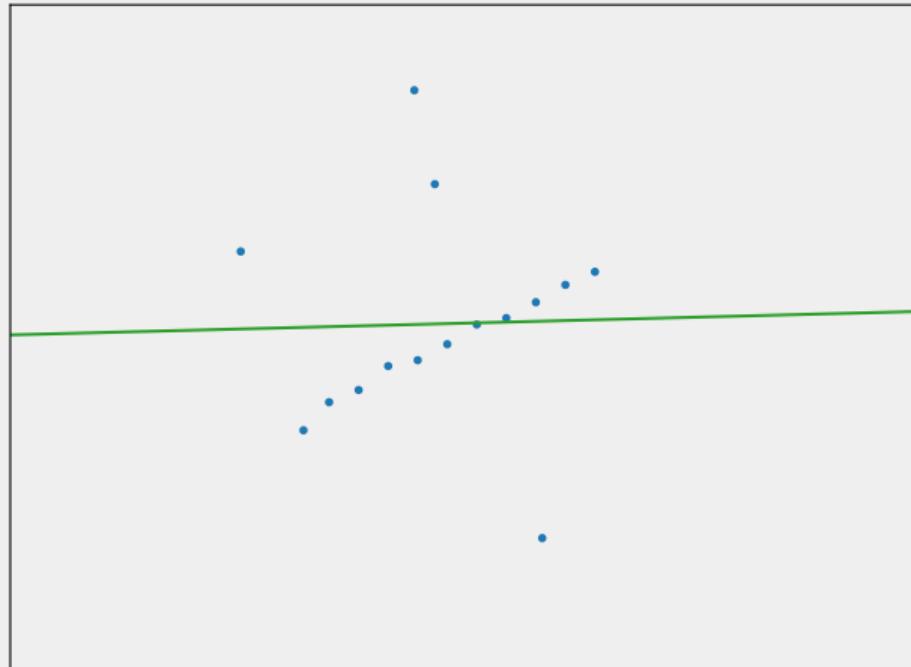


Fitting models

Can we fit a straight line
when there are a few
outliers?

Not really...

We get a very bad fit.
We need **robust** ways to fit
models!



About lines

- This presentation uses fitting straight lines to 2D points for all examples.
- Do we really care that much about fitting straight lines?

About lines

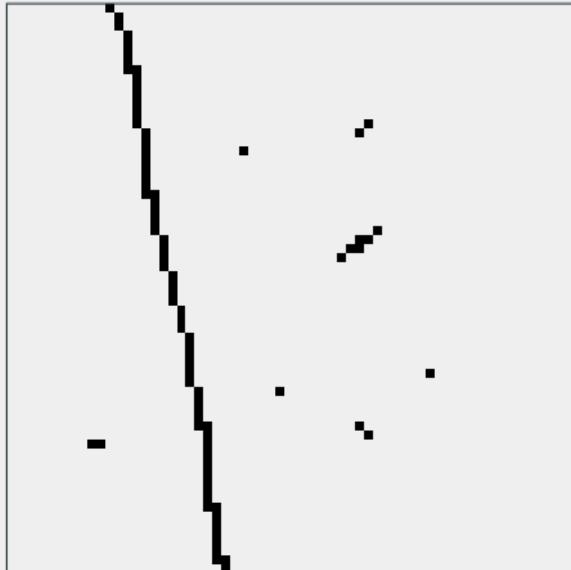
- This presentation uses fitting straight lines to 2D points for all examples.
- Do we really care that much about fitting straight lines?
- The principles *generalize* to other models!

Hough Transform

Hough Transform

Is a transformation of an edge image where lines can be extracted.

Example image



Hough transform of image



Hough Transform

How to represent a line?

- $y = ax + b$?
 - Has singularities for vertical lines
- Homogeneous coordinates?
 - Is over-parametrized

Hough Transform - r, θ representation

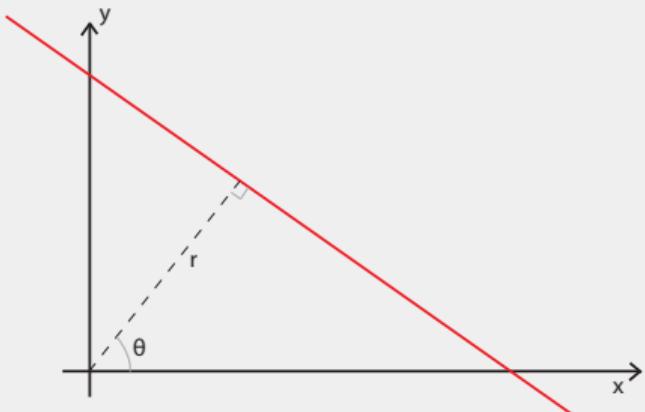
Represent using:

θ angle of the line

r closest distance from origin to line

Closely related to the homogeneous line representation

$$\mathbf{l} = [\cos(\theta) \quad \sin(\theta) \quad -r]$$



Hough Transform

We can represent any line with these two parameters.

- $\theta \in [-\frac{\pi}{2}, \frac{\pi}{2}[$
- $r \in [-d, d]$, where d is the diagonal length of the image.

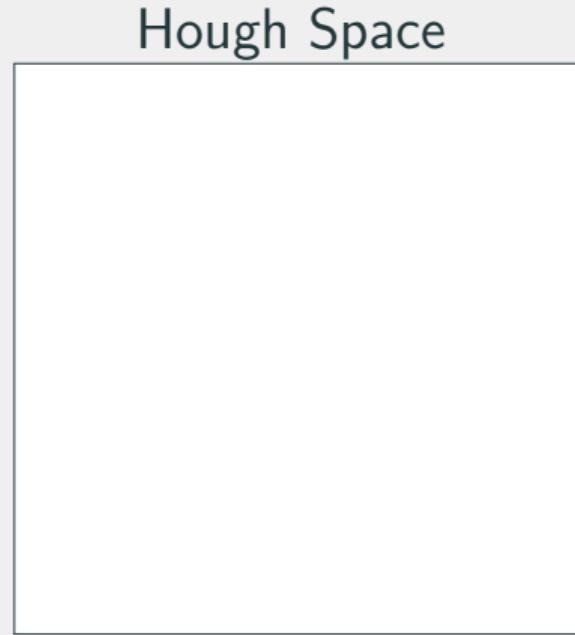
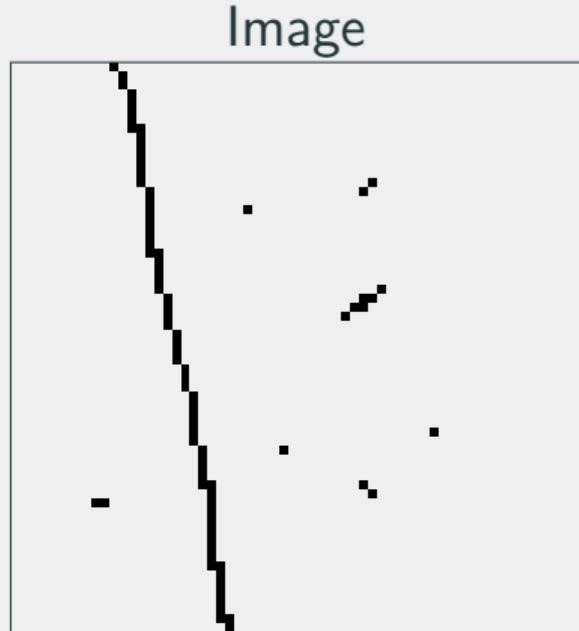
We discretize these values and represent them on a 2D grid

Hough Transform

Idea: Let points vote on which line is the best!

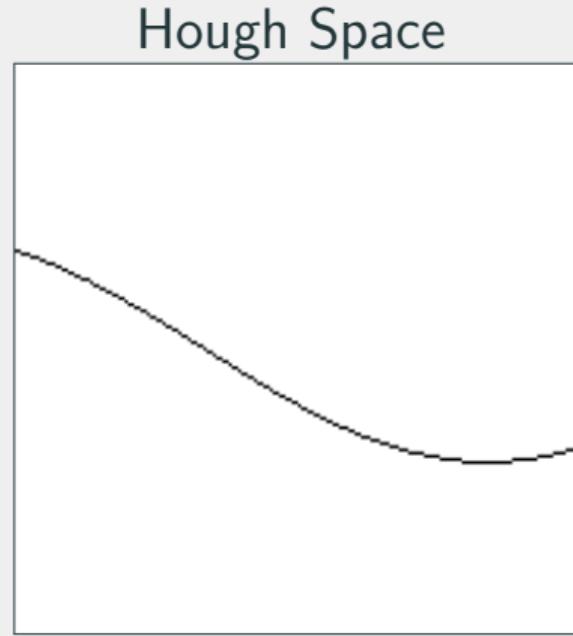
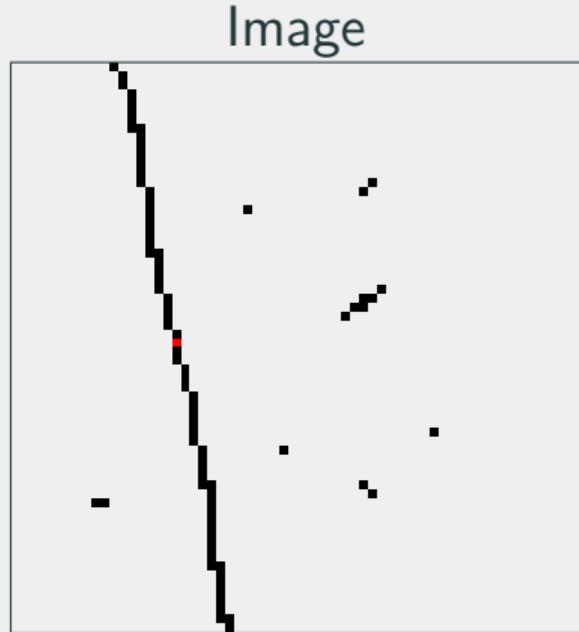
- Each point can be part of (infinitely) many lines.
- All potential lines going through this point are of interest
- Each point votes on all lines that go through its
 - This corresponds to a line in Hough space
- Repeat for all points

Hough Transform – Example

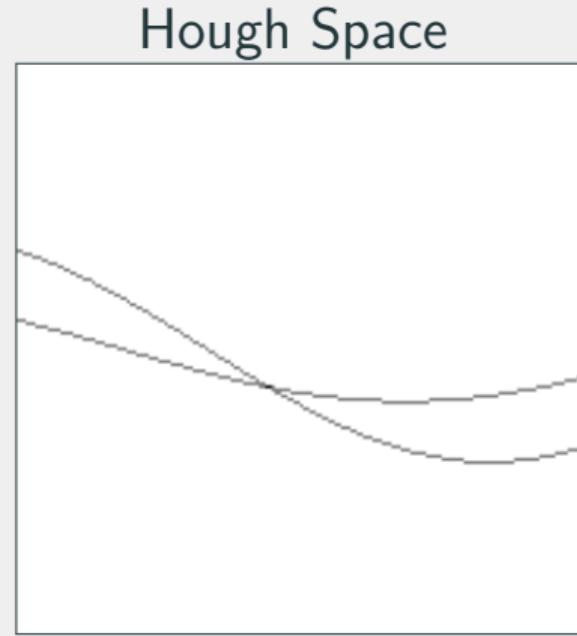
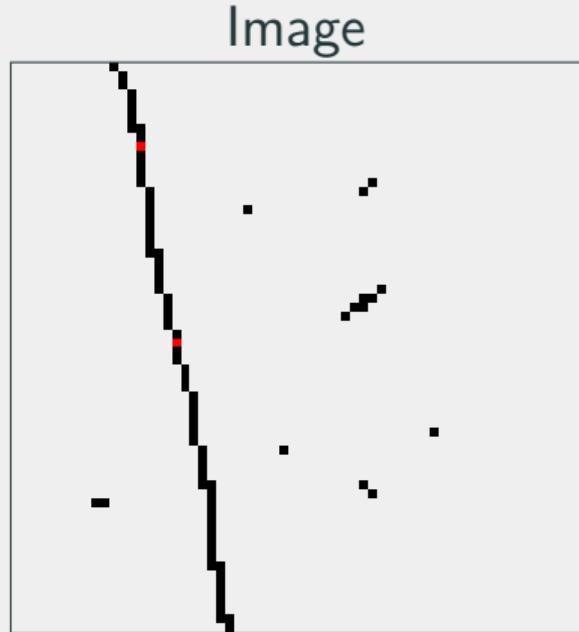


r

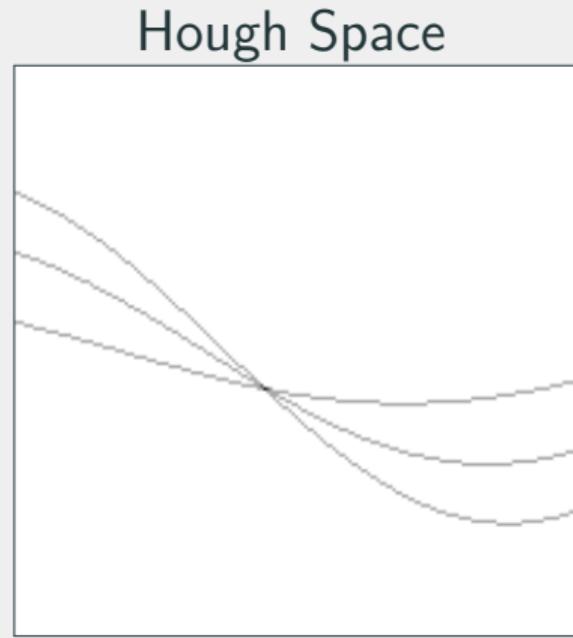
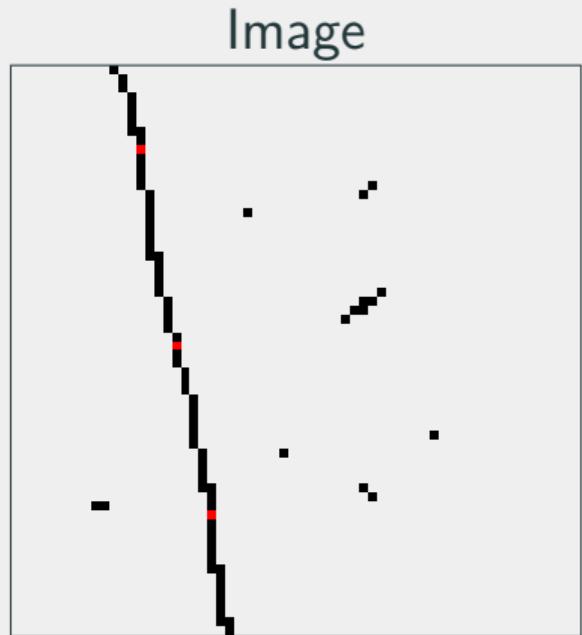
Hough Transform – Example



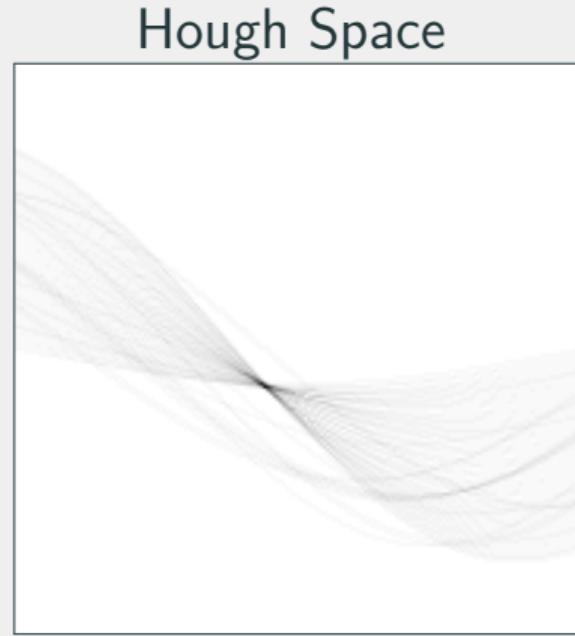
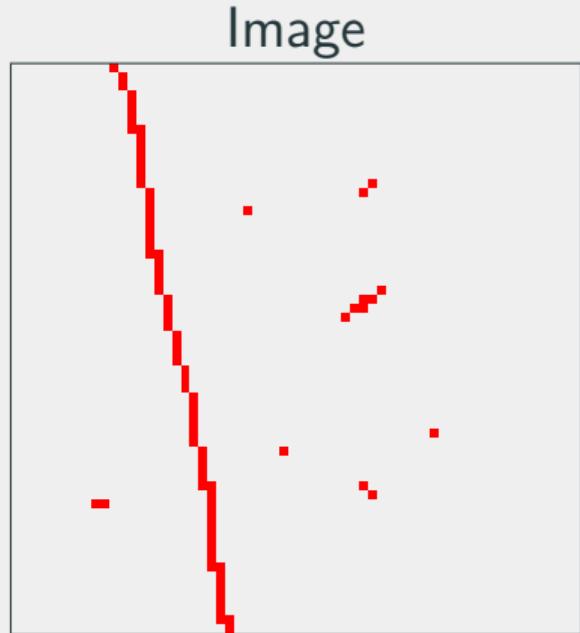
Hough Transform – Example



Hough Transform – Example



Hough Transform – Example



Hough Transform

- A peak in Hough space corresponds to a line in the image.
- Can be found using non-maximum suppression.

Generalized Hough Transform

- We can generalize the Hough transform for more complex models.
 - e.g. for circles, the Hough space is now in 3D and each point becomes a conic.
- Hough space has same number of dimensions as the model we fit has degrees of freedom.
- Impractical for more than three degrees of freedom

RANSAC

Random sample consensus (RANSAC)

Idea!

Instead of computing the Hough space, what if we could sample points directly in Hough space?

Random sample consensus (RANSAC)

Idea!

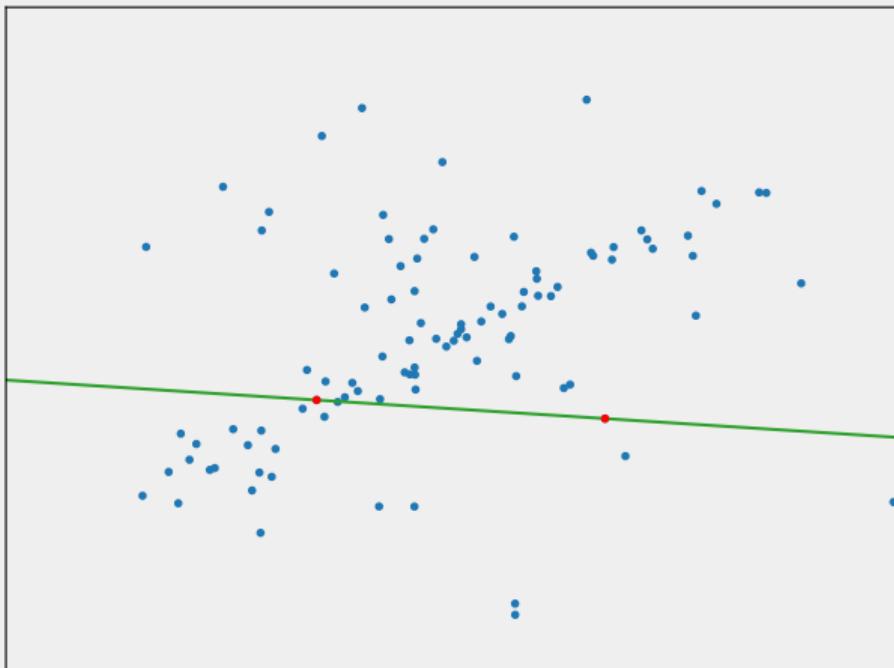
Instead of computing the Hough space, what if we could sample points directly in Hough space?

What if we could sample with the value in hough space being proportional to the probability of sampling the point?

RANSAC

- Randomly sample the minimum number of points we need to fit our model
- Fit the model to these samples

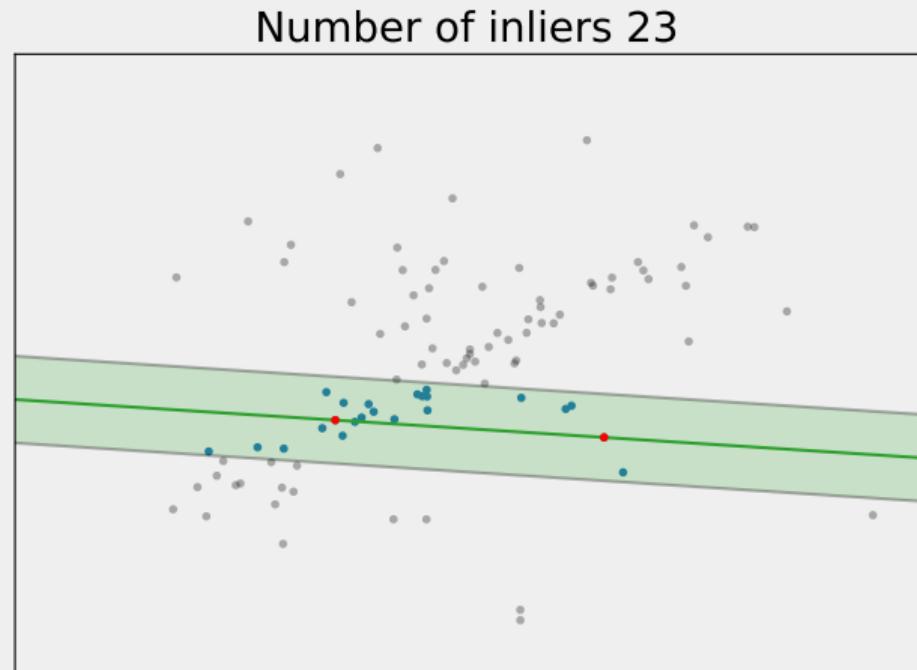
Does this line fit the data well?



Measure inliers

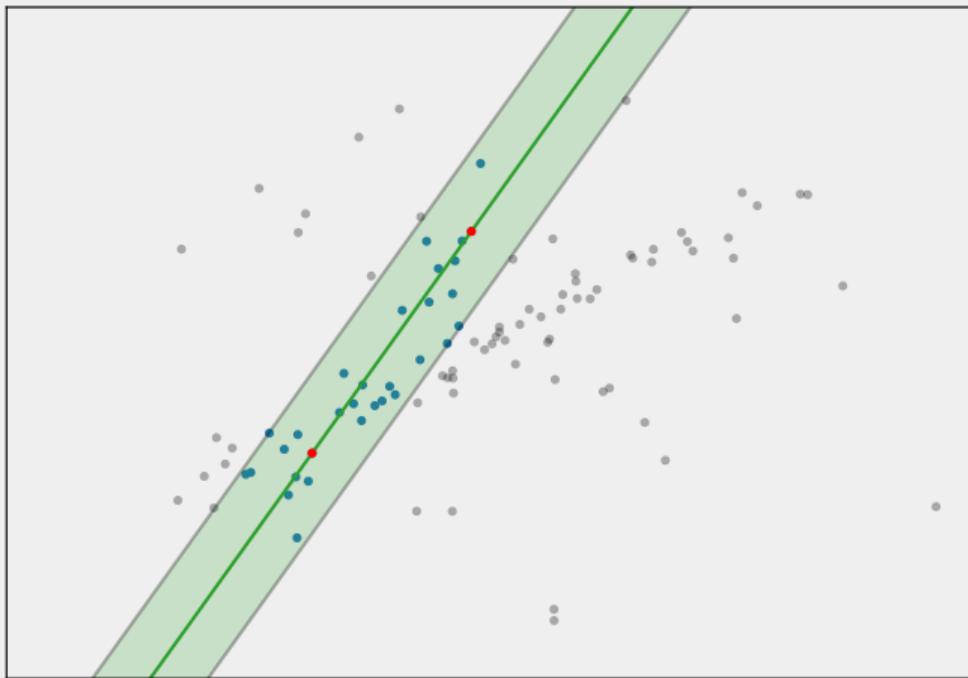
Points closer than a certain threshold to the line are **inliers**!

Number of inliers is indication of how well the line fits.



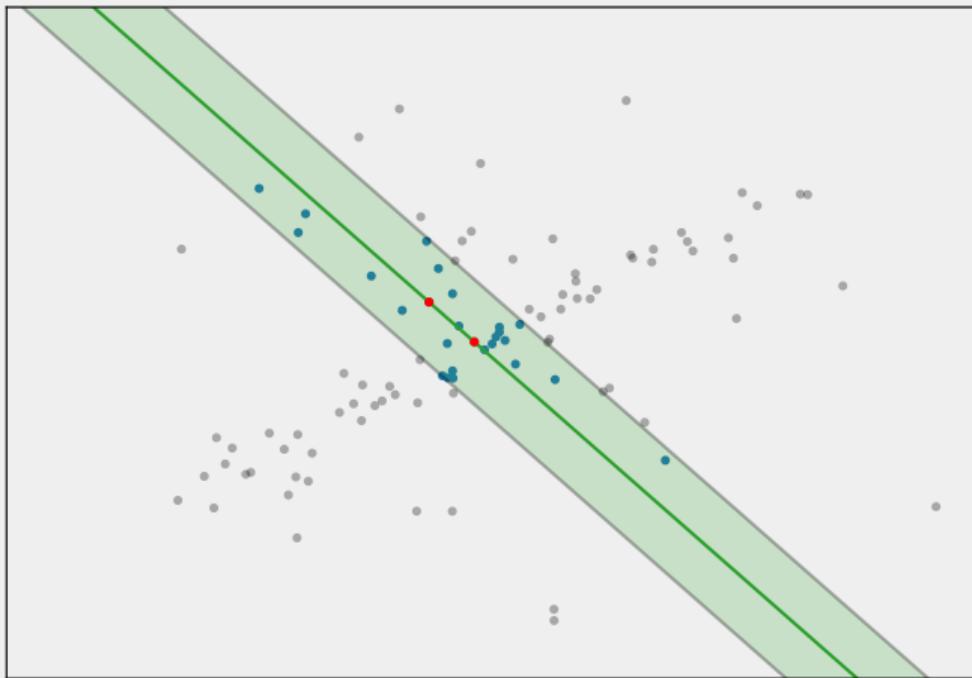
Example line

Number of inliers 32



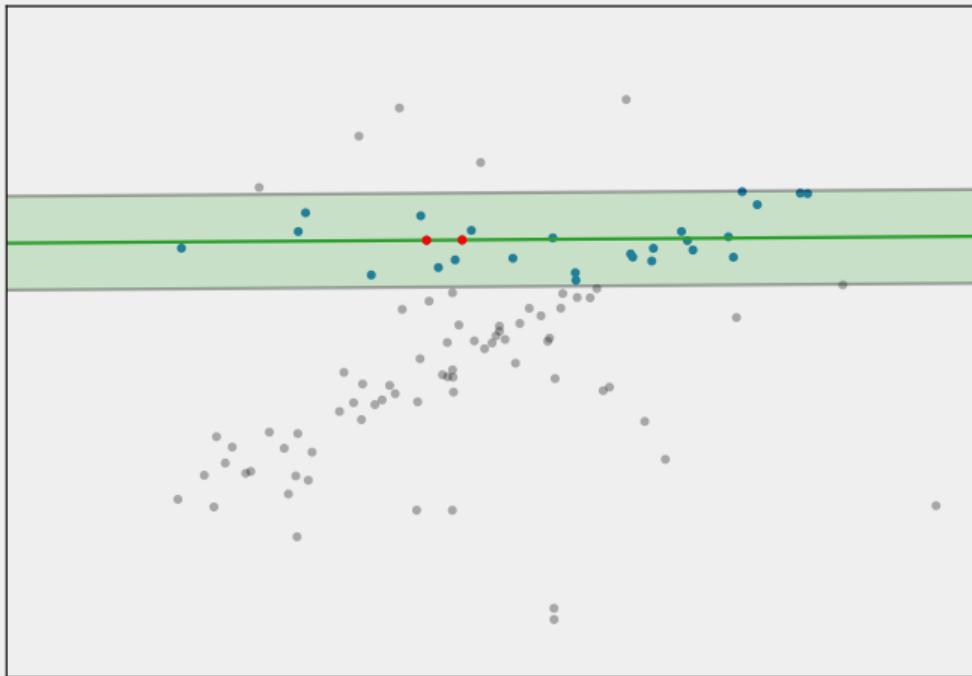
Example line

Number of inliers 26



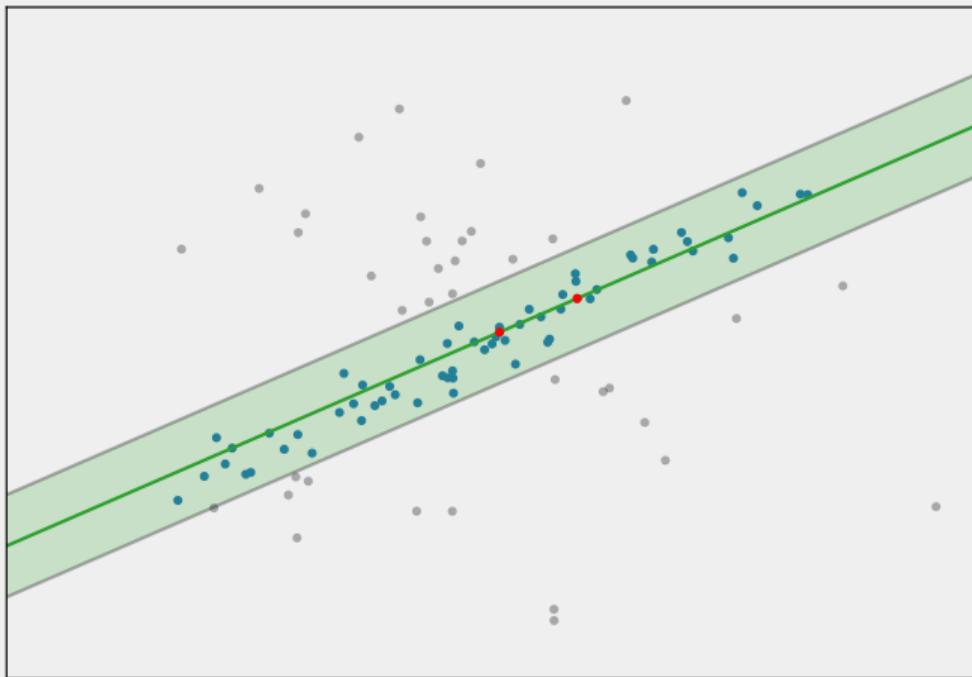
Example line

Number of inliers 27



Example line

Number of inliers 62



The RANSAC algorithm

Keep track of which fit had the most inliers so far, and update this if a fit with more inliers is found

The RANSAC algorithm

- Sample minimum number of points required to fit model
 - Fit model to these
- Data points with an error less than τ are inliers with respect to fitted model
 - If number of inliers is higher than the highest number of inliers seen so far, update best model.

The RANSAC algorithm

- Sample minimum number of points required to fit model
 - Fit model to these
- Data points with an error less than τ are inliers with respect to fitted model
 - If number of inliers is higher than the highest number of inliers seen so far, update best model.
- Repeat for N iterations.
- Final step:
 - Re-fit model to all inliers of the best model

Implementation details

- Represent the lines using homogeneous coordinates $[a \ b \ c]$
 - Makes it easy to compute distance to line
 - Scale the line such that $a^2 + b^2 = 1$
- Recall that distance to line is given by: $|\mathbf{l} \cdot \Pi^{-1}(\mathbf{p})|$
 - When $a^2 + b^2 = 1$

RANSAC

- Sample in Hough space without computing it.
- Useful for fitting models when outliers are present.
- We must select the threshold for inliers and the number of iterations carefully.
- Being able to fit a model to the least amount of data points is of interest

How many iterations?

- We can come with some idea of how many iterations we need
- Assume fraction of outliers is ϵ
- i.e. $\epsilon = 0.1$ means 10% of data are outliers.

How many iterations?

- We need n data points to fit a single model

$$P(\text{one sample has only inliers}) = (1 - \epsilon)^n$$

- Set p to a high value such as $p = 0.99$. Useful if we know ϵ .

How many iterations?

- We need n data points to fit a single model

$$P(\text{one sample has only inliers}) = (1 - \epsilon)^n$$

$$P(\text{one sample has one or more outliers}) = 1 - (1 - \epsilon)^n$$

- Set p to a high value such as $p = 0.99$. Useful if we know ϵ .

How many iterations?

- We need n data points to fit a single model

$$P(\text{one sample has only inliers}) = (1 - \epsilon)^n$$

$$P(\text{one sample has one or more outliers}) = 1 - (1 - \epsilon)^n$$

$$P(N \text{ samples all have one or more outliers}) = (1 - (1 - \epsilon)^n)^N$$

- Set p to a high value such as $p = 0.99$. Useful if we know ϵ .

How many iterations?

- We need n data points to fit a single model

$$P(\text{one sample has only inliers}) = (1 - \epsilon)^n$$

$$P(\text{one sample has one or more outliers}) = 1 - (1 - \epsilon)^n$$

$$P(N \text{ samples all have one or more outliers}) = (1 - (1 - \epsilon)^n)^N$$

$$P(\text{at least one of } N \text{ samples has only inliers}) = 1 - (1 - (1 - \epsilon)^n)^N$$

- Set p to a high value such as $p = 0.99$. Useful if we know ϵ .

How many iterations?

- We need n data points to fit a single model

$$P(\text{one sample has only inliers}) = (1 - \epsilon)^n$$

$$P(\text{one sample has one or more outliers}) = 1 - (1 - \epsilon)^n$$

$$P(N \text{ samples all have one or more outliers}) = (1 - (1 - \epsilon)^n)^N$$

$$P(\text{at least one of } N \text{ samples has only inliers}) = 1 - (1 - (1 - \epsilon)^n)^N = p$$

$$\Leftrightarrow N = \frac{\log(1 - p)}{\log(1 - (1 - \epsilon)^n)}$$

- Set p to a high value such as $p = 0.99$. Useful if we know ϵ .

Determining number of iterations adaptively

- We can estimate an upper bound of ϵ while running RANSAC.
- Let be s the number of inliers in the best model found so far
- Let be m the total number of data points
- $\hat{\epsilon} = 1 - \frac{s}{m} > \epsilon$

Determining number of iterations adaptively

- We can estimate an upper bound of ϵ while running RANSAC.
- Let be s the number of inliers in the best model found so far
- Let be m the total number of data points
- $\hat{\epsilon} = 1 - \frac{s}{m} > \epsilon$
- We can now estimate an upper bound on the number of iterations required
- $\hat{N} = \frac{\log(1-p)}{\log(1-(1-\hat{\epsilon})^n)} > N$
- Terminate once we have done more than \hat{N} iterations.

Recap of lines in homogeneous coordinates

- How do we fit a line to two points (p_1 and p_2) in homogeneous coordinates?
- What do we know about the line l ?
- $l \cdot p_1 = 0$ and $l \cdot p_2 = 0$.

Recap of lines in homogeneous coordinates

- How do we fit a line to two points (p_1 and p_2) in homogeneous coordinates?
- What do we know about the line l ?
- $l \cdot p_1 = 0$ and $l \cdot p_2 = 0$.
- The dot product is zero only when two vectors are perpendicular.
- How can we find l such that it is perpendicular to p_1 and p_2 ?

Recap of lines in homogeneous coordinates

- How do we fit a line to two points (p_1 and p_2) in homogeneous coordinates?
- What do we know about the line l ?
- $l \cdot p_1 = 0$ and $l \cdot p_2 = 0$.
- The dot product is zero only when two vectors are perpendicular.
- How can we find l such that it is perpendicular to p_1 and p_2 ?
- Use the **cross product**! $l = p_1 \times p_2$.

Learning objectives

After this lecture you should be able to:

- explain how the Hough transform works
- understand and implement RANSAC

Midterm Evaluation

Exercise time!

BLOBs and SIFT features

Morten R. Hannemose, mohan@dtu.dk

March 28, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- implement and use BLOB detection using Difference-of-Gaussians
- analyse and use SIFT features and feature matching

Similarity

Basic idea

- Locally appearance between views is the same
- Variation can be handled via invariances



Local image features

SIFT – key elements

- Features localized at interest points
- Adapted to scale and invariant to appearance changes



SIFT – scale invariant feature transform (Lowe, 1999)

- Scale-space BLOB detection – difference of Gaussians
- Interest point localization
- Orientation assignment
- Interest point descriptor
- Note – SIFT is one example of interest point feature

Harris corners and BLOBs

Harris corners are features that have a large change of intensity in two orthogonal directions. They are:

- local,
- can be found at different scales by changing the Gaussian filters, and
- invariant to rotation.

Harris corners are found by first order derivatives whereas BLOBs are response to second order image derivatives.

BLOBs – Binary Large OBjects

Correspond to:

- a dark area surrounded by brighter intensities or,
- a bright area surrounded by darker intensities

Hessian

The Hessian matrix contains the second order derivatives

$$\mathbf{H}(x, y) = \begin{bmatrix} I_{xx}(x, y) & I_{xy}(x, y) \\ I_{xy}(x, y) & I_{yy}(x, y) \end{bmatrix},$$

where

$$I_{xx}(x, y) = \frac{\partial^2 I(x, y)}{\partial x^2}, \quad I_{yy}(x, y) = \frac{\partial^2 I(x, y)}{\partial y^2}, \quad \text{and} \quad I_{xy}(x, y) = \frac{\partial^2 I(x, y)}{\partial x \partial y}.$$

Curvature

Second order derivatives measure **curvature**.

Eigenvalues of the Hessian (λ_1, λ_2) measure the principal curvature, i.e. the degree of change in derivative.

The eigenvectors measure the direction of that change

- the eigenvector corresponding to the largest eigenvalue (λ_1) is the direction of most change
- the second is orthogonal to that.

BLOB detection with Hessian

Similar to the Harris corner detector, we can use either of the measures

$$\det(\mathbf{H}) = \lambda_1 \lambda_2,$$

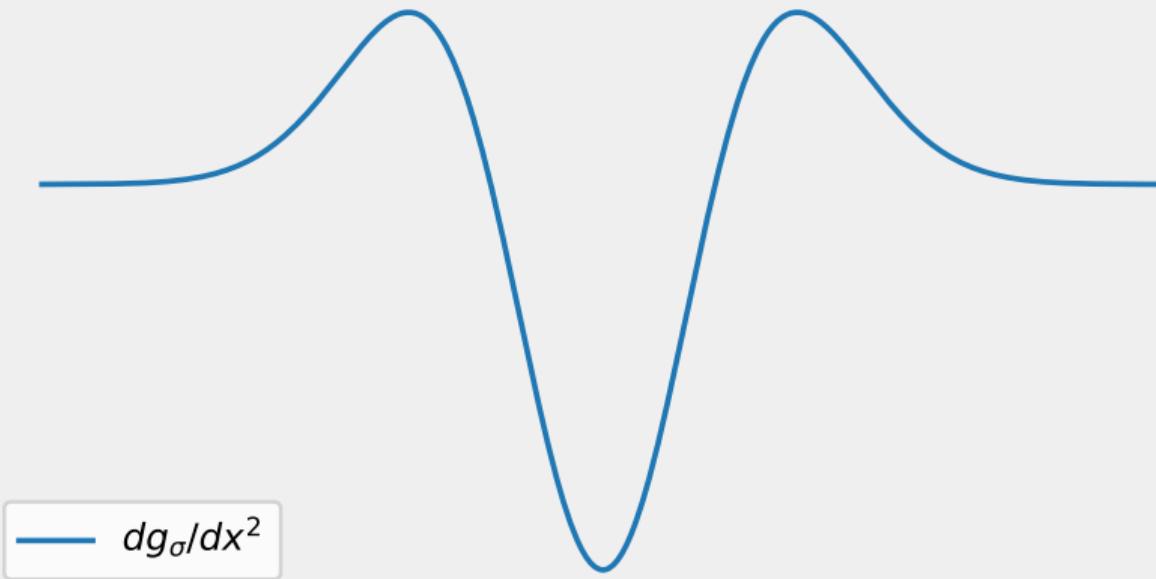
$$\text{trace}(\mathbf{H}) = \lambda_1 + \lambda_2,$$

where λ_i are the eigenvalues of the Hessian.

$\det(\mathbf{H})$ is the Gaussian curvature.

$\text{trace}(\mathbf{H}) = \nabla^2 I$ is **the Laplacian**, which we use for BLOB detection.

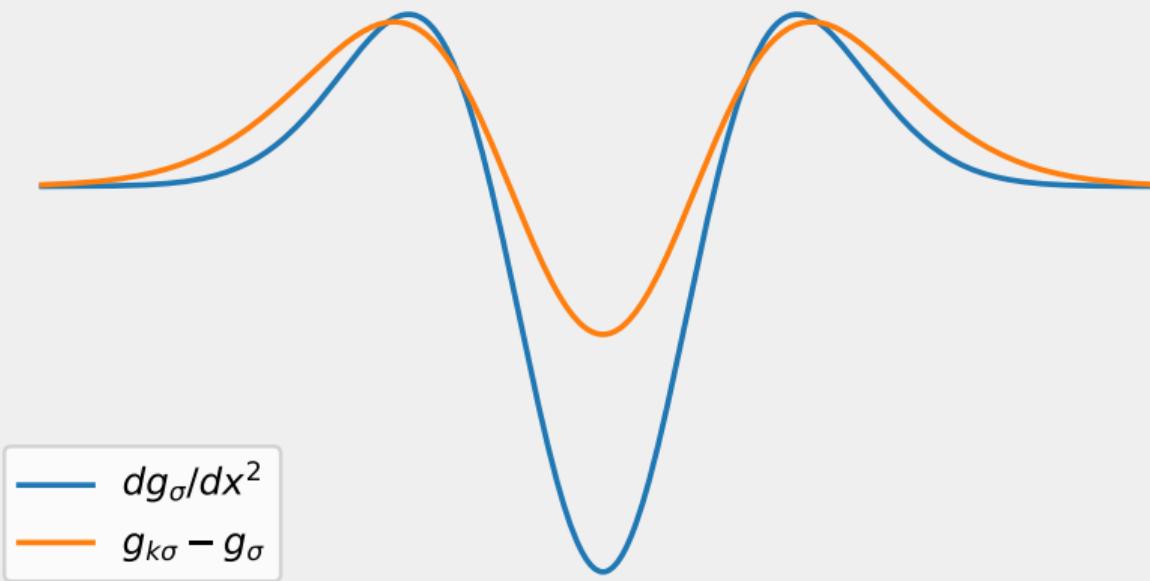
The Laplacian



Two Gaussians with different standard deviations



Difference of Gaussians vs Laplacian



BLOB detection with DoG

The Laplacian $\nabla^2 I$ can be approximated with the Difference-of-Gaussians (DoG).

Blurring the image with two different Gaussian kernels:

$$\nabla^2 I \approx D_\sigma = (G_{k\sigma} - G_\sigma) * I = G_{k\sigma} * I - G_\sigma * I,$$

where G_σ is a Gaussian with standard deviation σ and $k > 1$ is a scale factor.

BLOB detection with DoG

The Laplacian $\nabla^2 I$ can be approximated with the Difference-of-Gaussians (DoG).

Blurring the image with two different Gaussian kernels:

$$\nabla^2 I \approx D_\sigma = (G_{k\sigma} - G_\sigma) * I = G_{k\sigma} * I - G_\sigma * I,$$

where G_σ is a Gaussian with standard deviation σ and $k > 1$ is a scale factor.

Why are we interested in this approximation?

BLOB detection with DoG



Figure 1: A dog detecting blobs

SIFT – Scale invariance



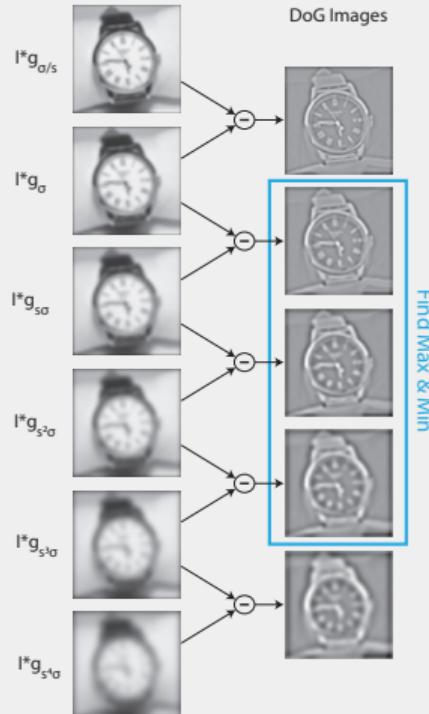
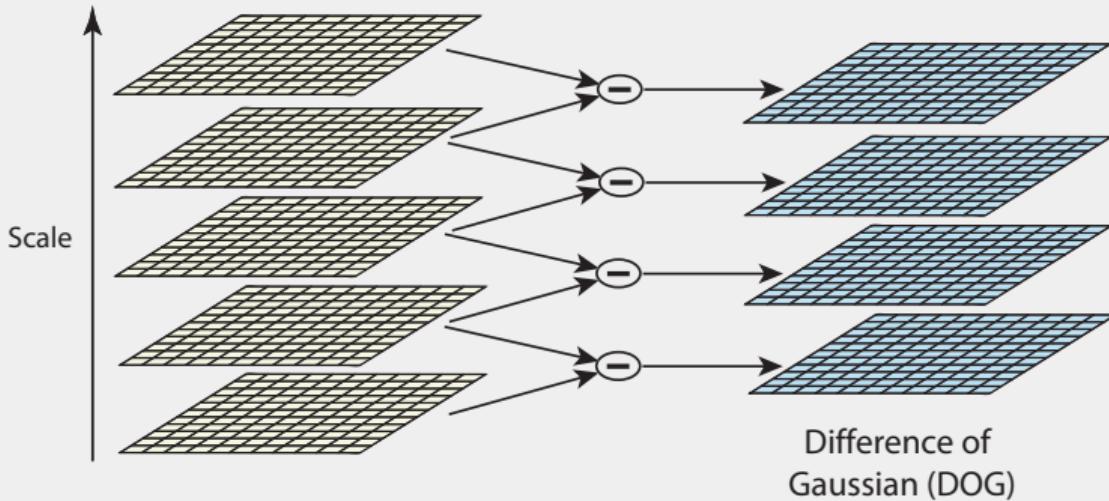
SIFT – Scale invariance

- Using difference of Gaussians for BLOB detection

$$\begin{aligned} D(x, y, \sigma) &= ((G_{k\sigma} - G_\sigma) * I)(x, y) \\ &= (G_{k\sigma} * I)(x, y) - (G_\sigma * I)(x, y) \\ &= L(x, y, k\sigma) - L(x, y, \sigma) \end{aligned}$$

The DoG is computed by subtracting more and more blurred images from each other.

SIFT - Difference of Gaussians



Gaussian scale space – Efficient

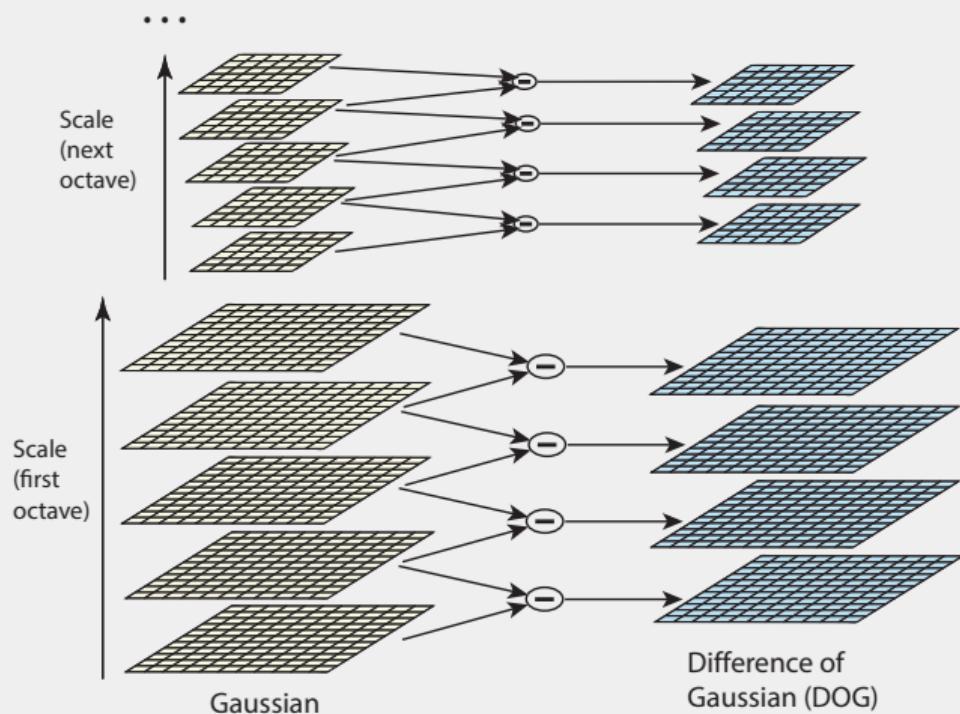
- Convolution of two Gaussians yield a new Gaussian
- Generate scale space by iteratively blurring already blurred images again
 - Otherwise we would need very large Gaussian kernels
- However the size of the kernel still grows
 - $(k\sigma)^2 = \sigma^2 + new^2$

Gaussian scale space – Efficient

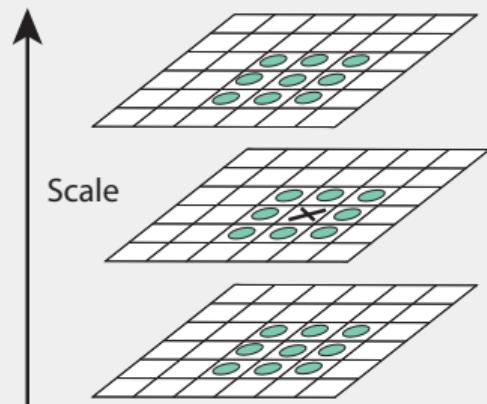
- Convolution of two Gaussians yield a new Gaussian
- Generate scale space by iteratively blurring already blurred images again
 - Otherwise we would need very large Gaussian kernels
- However the size of the kernel still grows
 - $(k\sigma)^2 = \sigma^2 + new^2$
- We choose $k = 2^{\frac{1}{3}}$
- σ doubles after three images, the image is downsampled.
 - This is an **octave**.
 - We only need three Gaussians of constant size (precomputed)

SIFT – Estimation of DoG

Difference of Gaussians



Extrema localization



SIFT - Magnitude of the DoG response

Do we get smaller values in the difference of Gaussians for high values of σ ?

Using the heat equation it can be shown that the response does not change as a function of sigma.

We can use the same threshold for the entire scale space.

SIFT – Subpixel localization

- Why is this necessary?

SIFT – Subpixel localization

- Why is this necessary?
- Second order Taylor approximation of DoG around local maximum:

$$D(\mathbf{x}) = D + \frac{\partial D^T}{\partial \mathbf{x}} \mathbf{x} + \frac{1}{2} \mathbf{x}^T \frac{\partial^2 D}{\partial \mathbf{x}^2} \mathbf{x}$$

- Setting the derivative of $D(\mathbf{x})$ to zero

$$\hat{\mathbf{x}} = -\frac{\partial^2 D^{-1}}{\partial \mathbf{x}^2} \frac{\partial D}{\partial \mathbf{x}}$$

SIFT – Subpixel localization

- We get

$$D(\hat{\mathbf{x}}) = D + \frac{1}{2} \frac{\partial D^T}{\partial \mathbf{x}} \hat{\mathbf{x}}$$

- If $|D(\hat{\mathbf{x}})| < 0.03$ the point is discarded
 - Removes points with low contrast.

SIFT – Interest point along edges discarded

- The eigenvalues of the Hessian are proportional to the principal curvatures

$$\mathbf{H} = \begin{bmatrix} I_{xx} & I_{xy} \\ I_{xy} & I_{yy} \end{bmatrix} = \begin{bmatrix} \frac{\partial^2 I}{\partial x^2} & \frac{\partial^2 I}{\partial x \partial y} \\ \frac{\partial^2 I}{\partial x \partial y} & \frac{\partial^2 I}{\partial y^2} \end{bmatrix}$$

Interest point along edges discarded

- λ_1 and λ_2 are the eigenvalues of the Hessian

$$\text{trace}(\mathbf{H}) = I_{xx} + I_{yy} = \alpha + \beta$$

$$\det(\mathbf{H}) = I_{xx}I_{yy} - I_{xy}^2 = \alpha\beta$$

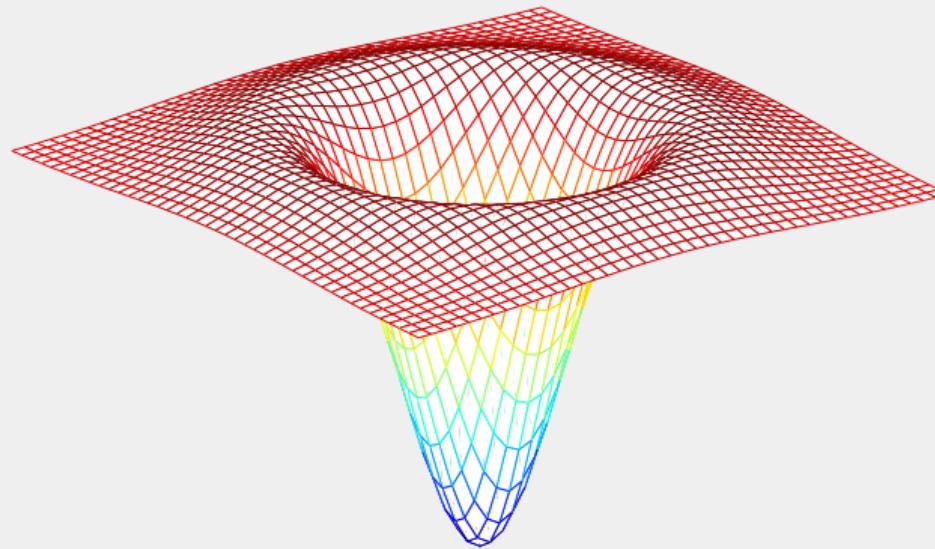
Points are kept if

$$\frac{\text{trace}(\mathbf{H})^2}{\det(\mathbf{H})} < \frac{(r+1)^2}{r},$$

where $r = 10$ (found to be a good heuristic)

Short break

DoG measure

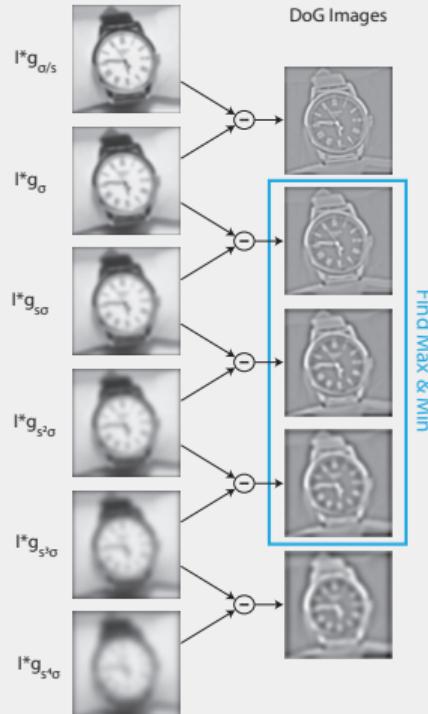


Features have $|D(\mathbf{x})| > \tau$, where τ is a threshold.
Dark BLOBs: $D(\mathbf{x}) > 0$, bright BLOBs: $D(\mathbf{x}) < 0$

DoG in scale pyramids

Scale pyramids are increasingly blurred of the same image.

Scale space DoG is subtraction between all layers adjacent scales.



Scale space BLOBs and DoG

DoGs at different scales makes for a scale invariant feature detector.

Small and large details are recoverable in different DoGs.



SIFT – Orientation assignment

Compute the orientation of gradients in a small region around the BLOB.

$$m(x, y) = \sqrt{L_x^2 + L_y^2}$$

$$\theta(x, y) = \arctan 2(L_y, L_x)$$

Where

$$L_x = L(x + 1, y) - L(x - 1, y)$$

$$L_y = L(x, y + 1) - L(x, y - 1)$$

SIFT – Orientation assignment

- Compute circular histogram of gradient orientations
 - Weighted by magnitude, smoothed, and has 36 bins

SIFT – Orientation assignment

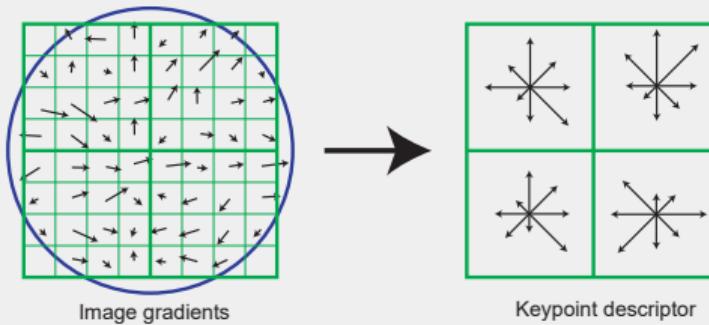
- Compute circular histogram of gradient orientations
 - Weighted by magnitude, smoothed, and has 36 bins
- Use peak in histogram to assign orientation of point
- This introduces **rotation invariance**.
- Can we have multiple peaks in histogram?

SIFT – Orientation assignment

- Compute circular histogram of gradient orientations
 - Weighted by magnitude, smoothed, and has 36 bins
- Use peak in histogram to assign orientation of point
- This introduces **rotation invariance**.
- Can we have multiple peaks in histogram?
 - Yes, this can happen at e.g. corners.
 - Create a new point at the same location if peak is over 80% of max.

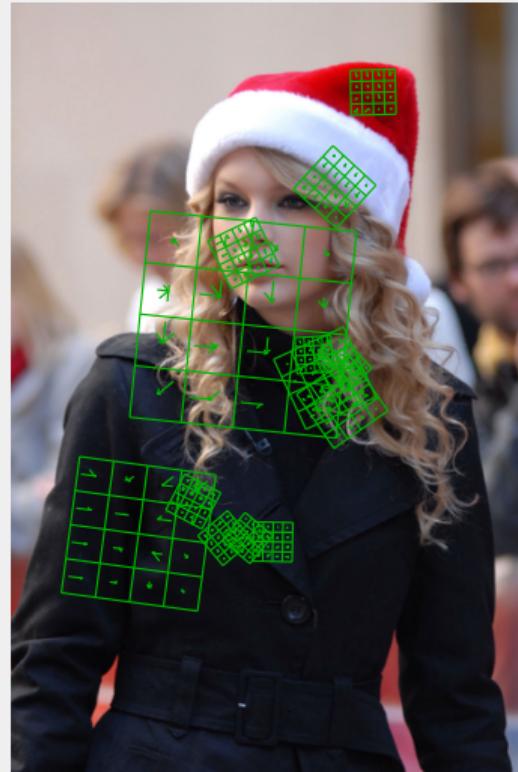
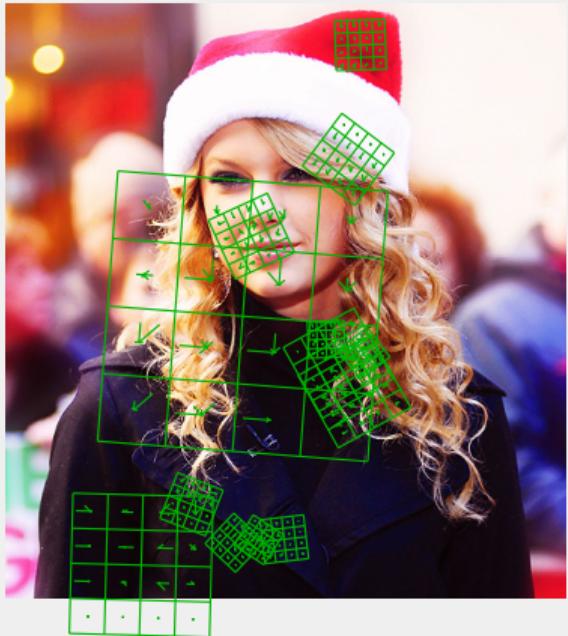
SIFT – Descriptor

- Create local patch at scale and orientation of point
- Build a histogram of local gradient orientations

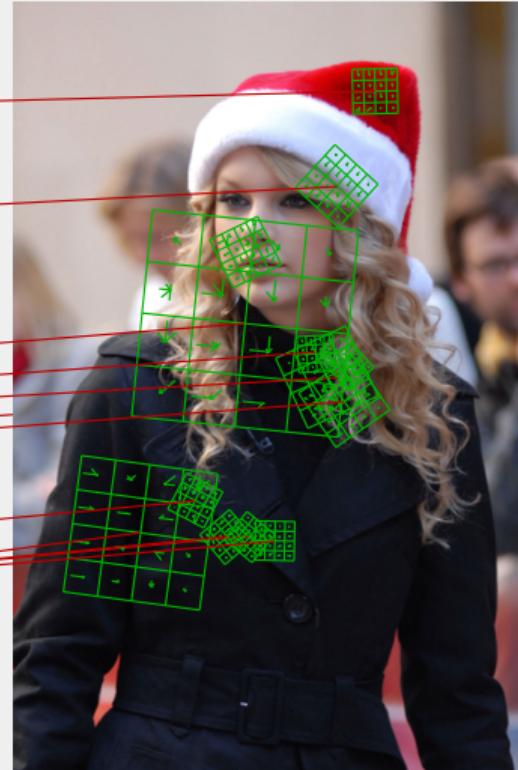
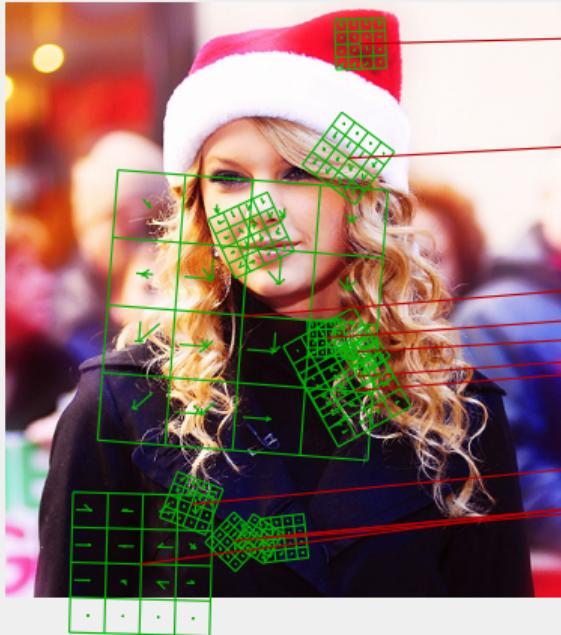


- Normalized using L₂ norm: $\mathbf{d}_n = \frac{1}{\sqrt{\sum_{i=1}^{128} \mathbf{d}(i)^2}} \mathbf{d}$

Taylor SIFT



Taylor SIFT



SIFT – Invariances

- Position
- Scale
- Rotation
- Linear intensity change
- Perspective changes?

SIFT – Matching of descriptors

- Use Euclidean distance between normalized vectors

$$\delta(\mathbf{d}_i, \mathbf{d}_j) = \sqrt{\sum_{n=1}^{128} (d_{i,n} - d_{j,n})^2}$$

- Note – for comparison the square root is not needed

RootSIFT

Simple trick to improve SIFT matching

- SIFT is a histogram
- Euclidean distance is dominated by large values
- RootSIFT is a transformation that measures distance using the Hellinger kernel.
 - L1 normalize
 - Take the square root of each element
 - L2 normalize the resulting vector
- Compare using Euclidean distance

SIFT – Matching of descriptors

- For each feature in image 1 ($d_{1,i}$) find the closest feature in image 2 ($d_{2,j}$)
 - This will give a lot of incorrect matches
- Cross checking
 - Only keep matches where $d_{2,j}$ is also the closest to $d_{1,i}$ of all features in image 1

SIFT – Matching of descriptors

- For each feature in image 1 ($d_{1,i}$) find the closest feature in image 2 ($d_{2,j}$)
 - This will give a lot of incorrect matches
- Cross checking
 - Only keep matches where $d_{2,j}$ is also the closest to $d_{1,i}$ of all features in image 1
- Ratio test
 - Compute the ratio between the closest and second closest match, and keep where this is below a threshold, e.g. 0.7.

SIFT – Summary

- SIFT is both a feature detector and descriptor
- Find local extrema of DoGs in scale space
- Place patch oriented along local gradients
- Compute histograms of gradients.
- Allows matching of images invariant to: scale, rotation, illumination and viewpoint
- Partly visible objects can be matched

Other descriptors

- SIFT is widely used. (74k+ citations)
 - Was patented until 2020.
- Meanwhile other similar methods were created
 - SURF, 2008 (14k+ citations)
 - ORB 2011, (12k+ citations)
 - BRIEF 2010, (5k+ citations)
 - BRISK 2011, (4k+ citations)

Learned descriptors

- Deep Learning has created improved feature detectors/descriptors.
- Mostly in improvement in invariance to changing lighting.
- Some examples:
 - R2D2: Repeatable and Reliable Detector and Descriptor [\[code\]](#)
 - Superpoint [\[code\]](#)

Learning objectives

After this lecture you should be able to:

- implement and use BLOB detection using Difference-of-Gaussians
- analyse and use SIFT features and feature matching

Exercise

Build a BLOB detector and match points with SIFT detector.

Python: Use OpenCV (4.2.0 or newer)

Matlab: Use VLFeat

<https://www.vlfeat.org/overview/sift.html>

Exercise time!

Geometry Constrained Feature Matching

Morten R. Hannemose, mohan@dtu.dk

April 4, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- explain and implement the eight point algorithm for estimating the fundamental matrix
- explain and implement estimation of the fundamental matrix with RANSAC
- choose the threshold for RANSAC using χ^2

Presentation topics

Estimating the fundamental matrix

Linear algorithm

Incorporating RANSAC

Thresholding for RANSAC

Setting the scene

- Stereo geometry
- (SIFT) Features
- RANSAC

Geometry constrained feature matching

Use multi view geometry to filter matches



Fundamental matrix – Recap

R and t describe the relative pose between the cameras.

$$E = [t]_{\times} R$$

$$F = K_2^{-\top} E K_1^{-1}$$

$$0 = q_2^{\top} F q_1$$

The fundamental matrix expresses that two corresponding points lie on their corresponding epipolar lines.

Estimating the fundamental matrix

Fundamental matrix problem

The fundamental matrix is defined with the relation $\mathbf{q}_2^T \mathbf{F} \mathbf{q}_1 = 0$.

Consider the corresponding points \mathbf{q}_{1i} and \mathbf{q}_{2i} projected into cameras one and two, respectively. The relation is then

$$\begin{aligned} 0 &= \mathbf{q}_{2i}^T \mathbf{F} \mathbf{q}_{1i} \\ &= \begin{bmatrix} x_{2i} \\ y_{2i} \\ 1 \end{bmatrix}^T \begin{bmatrix} F_{11} & F_{12} & F_{13} \\ F_{21} & F_{22} & F_{23} \\ F_{31} & F_{32} & F_{33} \end{bmatrix} \begin{bmatrix} x_{1i} \\ y_{1i} \\ 1 \end{bmatrix} \end{aligned}$$

Fundamental matrix problem

Rearrange the terms:

$$0 = \mathbf{q}_{2i}^\top \mathbf{F} \mathbf{q}_{1i},$$

$$0 = \mathbf{B}^{(i)} \text{flatten}(\mathbf{F}),$$

where

$$\begin{aligned}\mathbf{B}^{(i)} &= [x_{1i}x_{2i} \quad y_{1i}x_{2i} \quad x_{2i} \quad x_{1i}y_{2i} \quad y_{1i}y_{2i} \quad y_{2i} \quad x_{1i} \quad y_{1i} \quad 1] \\ &= \text{flatten}(\mathbf{q}_{2i}\mathbf{q}_{1i}^\top),\end{aligned}$$

$$\text{flatten}(\mathbf{F}) = [F_{11} \quad F_{12} \quad F_{13} \quad F_{21} \quad F_{22} \quad F_{23} \quad F_{31} \quad F_{32} \quad F_{33}]^\top$$

Fundamental matrix solution

Define \mathbf{B}

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}^{(1)} \\ \mathbf{B}^{(2)} \\ \vdots \\ \mathbf{B}^{(n)} \end{bmatrix}.$$

Subject to $\|\text{flatten}(\mathbf{F})\|_2 = 1$ the solution is the singular vector with the smallest singular value.

Degrees of freedom – Eight point algorithm

F has 9 numbers, and is scale invariant.

Each pair of corresponding points fixes a degree of freedom. Eight points is enough to estimate the fundamental matrix.

This is the **eight point algorithm**.

Degrees of freedom – Seven point algorithm

$[t]_\times$ has rank 2, and thus \mathbf{F} is also rank deficient, i.e. $\det(\mathbf{F}) = 0$.

Thus \mathbf{F} has 7 degrees of freedom, and can be found from 7 matches.

\mathbf{B} will have two singular vectors with singular value 0.

Denote these \mathbf{F}' and \mathbf{F}^\dagger .

$\mathbf{F} = \alpha\mathbf{F}' + (1 - \alpha)\mathbf{F}^\dagger$, where α is chosen such that $\det(\mathbf{F}) = 0$.

This is the **seven point algorithm**.

Degrees of freedom – Seven point algorithm

$[t]_\times$ has rank 2, and thus \mathbf{F} is also rank deficient, i.e. $\det(\mathbf{F}) = 0$.

Thus \mathbf{F} has 7 degrees of freedom, and can be found from 7 matches.

\mathbf{B} will have two singular vectors with singular value 0.

Denote these \mathbf{F}' and \mathbf{F}^\dagger .

$\mathbf{F} = \alpha\mathbf{F}' + (1 - \alpha)\mathbf{F}^\dagger$, where α is chosen such that $\det(\mathbf{F}) = 0$.

This is the **seven point algorithm**. Discuss with your neighbour why it is preferable to estimate from 7 instead of 8.

F can be estimated from eight point correspondences easily.

Possible to estimate from just seven.

Estimating the fundamental matrix with RANSAC

To use RANSAC, we need a way to measure how well a datapoint fits a specific model. For lines, we measured the distance from a point to the line.

What is a datapoint when talking about the fundamental matrix?

Estimating the fundamental matrix with RANSAC

What does $q_{2i}^T F q_{1i}$ equal for non corresponding points?

Estimating the fundamental matrix with RANSAC

What does $\mathbf{q}_{2i}^T \mathbf{F} \mathbf{q}_{1i}$ equal for non corresponding points?

$\mathbf{q}_{2i}^T \mathbf{F}$ and $\mathbf{F} \mathbf{q}_{1i}$ are epipolar lines.

Distance from point to line

$$d = [a \ b \ c] [x \ y \ 1]^T$$

should have $a^2 + b^2 = 1$, this is not guaranteed for the epipolar lines.

Symmetric epipolar distance

We want to measure the distance

- from \mathbf{q}_{2i} to $\mathbf{F}\mathbf{q}_{1i}$ and
- from \mathbf{q}_{1i} to $\mathbf{F}^\top\mathbf{q}_{2i}$

We can compute these as

$$\frac{\mathbf{q}_{2i}^\top \mathbf{F} \mathbf{q}_{1i}}{\sqrt{(\mathbf{F} \mathbf{q}_{1i})_1^2 + (\mathbf{F} \mathbf{q}_{1i})_2^2}} \quad \text{and} \quad \frac{\mathbf{q}_{2i}^\top \mathbf{F} \mathbf{q}_{1i}}{\sqrt{(\mathbf{q}_{2i}^\top \mathbf{F})_1^2 + (\mathbf{q}_{2i}^\top \mathbf{F})_2^2}}$$

where \mathbf{x}_i^2 refers to the square of i^{th} element of \mathbf{x} .

Symmetric epipolar distance

We can normalize the distance to both epipolar lines, using their first two coordinates.

The squared symmetric epipolar distance is then given by

$$(\mathbf{q}_{2i}^\top \mathbf{F} \mathbf{q}_{1i})^2 \left(\frac{1}{(\mathbf{q}_{2i}^\top \mathbf{F})_1^2 + (\mathbf{q}_{2i}^\top \mathbf{F})_2^2} + \frac{1}{(\mathbf{F} \mathbf{q}_{1i})_1^2 + (\mathbf{F} \mathbf{q}_{1i})_2^2} \right),$$

Sampson's Distance

A similar distance is Sampson's distance.

$$d_{\text{Samp}}(\mathbf{F}, \mathbf{q}_{1i}, \mathbf{q}_{2i}) = \frac{(\mathbf{q}_{2i}^\top \mathbf{F} \mathbf{q}_{1i})^2}{(\mathbf{q}_{2i}^\top \mathbf{F})_1^2 + (\mathbf{q}_{2i}^\top \mathbf{F})_2^2 + (\mathbf{F} \mathbf{q}_{1i})_1^2 + (\mathbf{F} \mathbf{q}_{1i})_2^2}$$

Performs slightly better than the geometric distance in practice.

Is a squared distance.

Thresholding for RANSAC

Thresholding distances

- How to choose the threshold for RANSAC?

Thresholding distances

- How to choose the threshold for RANSAC?
- Introduce assumptions! 😊
- Assume that the errors of inliers follow a normal distribution.

Dimensionality of the error

- Fundamental/essential matrix
 - The error is the distance to the epipolar line
 - This is an error in one dimension
- Homography
 - The error is the distance from mapped point to true point
 - This is an error in two dimensions
- Pose estimation/camera calibration
 - Again it is a distance between points
 - This is an error in two dimensions

This is called the **codimension** of the problem.

We denote this m .

Choosing the threshold

- Assume that the error of each sample follows an m -dimensional normal distribution with standard deviation σ
- The summed and squared error is χ_m^2 distributed (by definition)
- Always work with squared distances
 - Not necessary to take square root when comparing distances

Choosing the threshold

- Assume that the error of each sample follows an m -dimensional normal distribution with standard deviation σ
- The summed and squared error is χ_m^2 distributed (by definition)
- Always work with squared distances
 - Not necessary to take square root when comparing distances
- Choose a confidence level e.g. 95%
 - i.e. we want our threshold to correctly identify 95% of true inliers.
- Look up the CDF for our χ_m^2 distribution.
- E.g. for a fundamental matrix and 95%, $\tau^2 = 3.84 \cdot \sigma^2$

Some values from the CDF of χ^2_m

$m \setminus 1 - p$	90%	95%	99%
1	2.71	3.84	6.63
2	4.61	5.99	9.21

Learning objectives

After this lecture you should be able to:

- explain and implement the eight point algorithm for estimating the fundamental matrix
- explain and implement estimation of the fundamental matrix with RANSAC
- choose the threshold for RANSAC using χ^2

Exercise time!

Image Stitching

Morten R. Hannemose, mohan@dtu.dk

April 11, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- explain and implement RANSAC to fit homographies
- understand explain the challenges involved in stitching panoramas

Presentation topics

Image stitching

Finding inliers

Transforming images

Multiple images

Nodal point

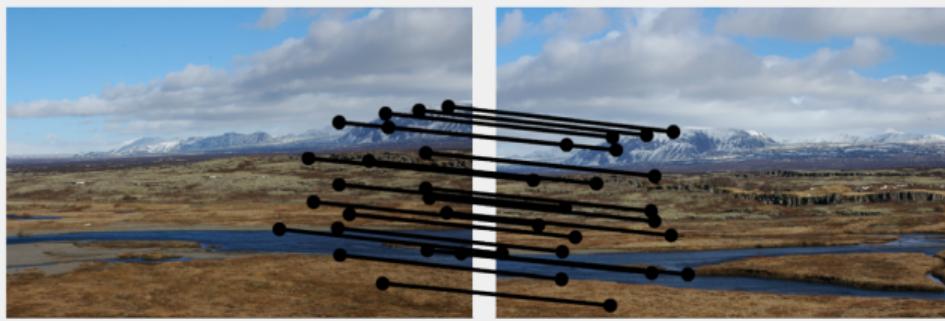
Setting the scene

Today you will be using:

- Homographies (hest)
- (SIFT) Features
- RANSAC

Image stitching

Image stitching – Panorama



About homographies

- What does a homography describe?

About homographies

- What does a homography describe?
 - The mapping of points between two images viewing the same plane.
- **BUT!** the mountains on the previous slide are not a plane!?
 - Why can we still use a homography?

Homographies for panoramas

- When the camera doesn't move but only rotates there are **no perspective deformations**
- It is equivalent to looking at a painting of the world
- Therefore we can assume the world is flat and use a homography!

Homographies for panoramas

- When the camera doesn't move but only rotates there are **no perspective deformations**
- It is equivalent to looking at a painting of the world
- Therefore we can assume the world is flat and use a homography!
- What happens if the camera moves?

Finding inliers

Measuring the error of a match

Let \mathbf{q}_1 and \mathbf{q}_2 be matching points without errors, such that

$$\mathbf{p}_1 = \mathbf{H}\mathbf{p}_2$$

$$\mathbf{p}_1 = \Pi^{-1}(\mathbf{q}_1), \quad \mathbf{p}_2 = \Pi^{-1}(\mathbf{q}_2)$$

$$\mathbf{q}_1 = \Pi(\mathbf{H}\Pi^{-1}(\mathbf{q}_2))$$

Measuring the error of a match

We have some 2D error in our detection of the points

$$\tilde{\mathbf{q}}_1 = \mathbf{q}_1 + \boldsymbol{\epsilon}_1, \quad \tilde{\mathbf{q}}_2 = \mathbf{q}_2 + \boldsymbol{\epsilon}_2$$

$$\tilde{\mathbf{p}}_1 = \Pi^{-1}(\tilde{\mathbf{q}}_1), \quad \tilde{\mathbf{p}}_2 = \Pi^{-1}(\tilde{\mathbf{q}}_2)$$

Measuring the error of a match

We have some 2D error in our detection of the points

$$\begin{aligned}\tilde{\mathbf{q}}_1 &= \mathbf{q}_1 + \boldsymbol{\epsilon}_1, & \tilde{\mathbf{q}}_2 &= \mathbf{q}_2 + \boldsymbol{\epsilon}_2 \\ \tilde{\mathbf{p}}_1 &= \Pi^{-1}(\tilde{\mathbf{q}}_1), & \tilde{\mathbf{p}}_2 &= \Pi^{-1}(\tilde{\mathbf{q}}_2)\end{aligned}$$

The error is the distance to the observed point in both images.

$$\|\tilde{\mathbf{q}}_1 - \Pi(\mathbf{H}\Pi^{-1}(\mathbf{q}_2))\|_2^2 + \|\tilde{\mathbf{q}}_2 - \mathbf{q}_2\|_2^2$$

Do we know \mathbf{q}_2 ?

Measuring the error of a match

We have some 2D error in our detection of the points

$$\begin{aligned}\tilde{\mathbf{q}}_1 &= \mathbf{q}_1 + \boldsymbol{\epsilon}_1, & \tilde{\mathbf{q}}_2 &= \mathbf{q}_2 + \boldsymbol{\epsilon}_2 \\ \tilde{\mathbf{p}}_1 &= \Pi^{-1}(\tilde{\mathbf{q}}_1), & \tilde{\mathbf{p}}_2 &= \Pi^{-1}(\tilde{\mathbf{q}}_2)\end{aligned}$$

The error is the distance to the observed point in both images.

$$\|\tilde{\mathbf{q}}_1 - \Pi(\mathbf{H}\Pi^{-1}(\mathbf{q}_2))\|_2^2 + \|\tilde{\mathbf{q}}_2 - \mathbf{q}_2\|_2^2$$

Do we know \mathbf{q}_2 ?

No!

Measuring the error of a match

We can use optimization to estimate the error free point

$$\text{dist}_{\text{true}}^2 = \min_{\tilde{\mathbf{q}}_2} \left\| \tilde{\mathbf{q}}_1 - \Pi(\mathbf{H}\Pi^{-1}(\mathbf{q}_2)) \right\|_2^2 + \left\| \tilde{\mathbf{q}}_2 - \mathbf{q}_2 \right\|_2^2$$

Requires solving a least squares problem for each distance computation.

Impractical!

Practical approximation

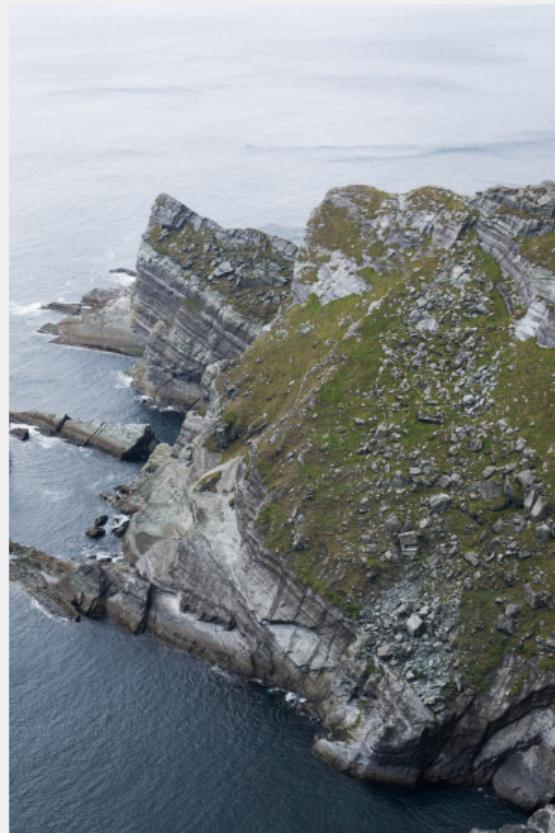
Map the observed points back and forth, and compute the distance

$$\text{dist}_{\text{approx}}^2 = \|\tilde{\mathbf{q}}_1 - \Pi(\mathbf{H}\tilde{\mathbf{p}}_2)\|_2^2 + \|\tilde{\mathbf{q}}_2 - \Pi(\mathbf{H}^{-1}\tilde{\mathbf{p}}_1)\|_2^2$$

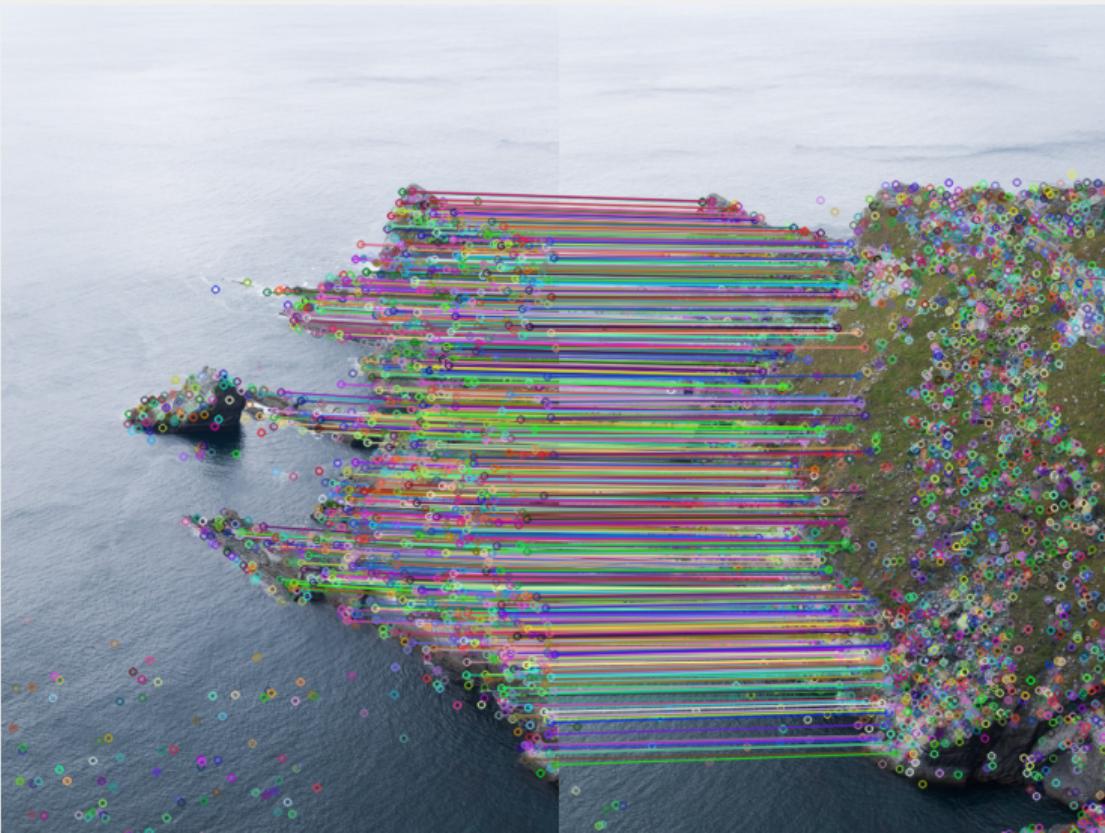
Use this distance to measure if a pair of points are inliers with respect to a homography.

Transforming images

Example – Images



Example – Matches



How to use the homography?

- We have estimated the homography
- How do we use it?
- Same approach as image undistortion.

Transforming images

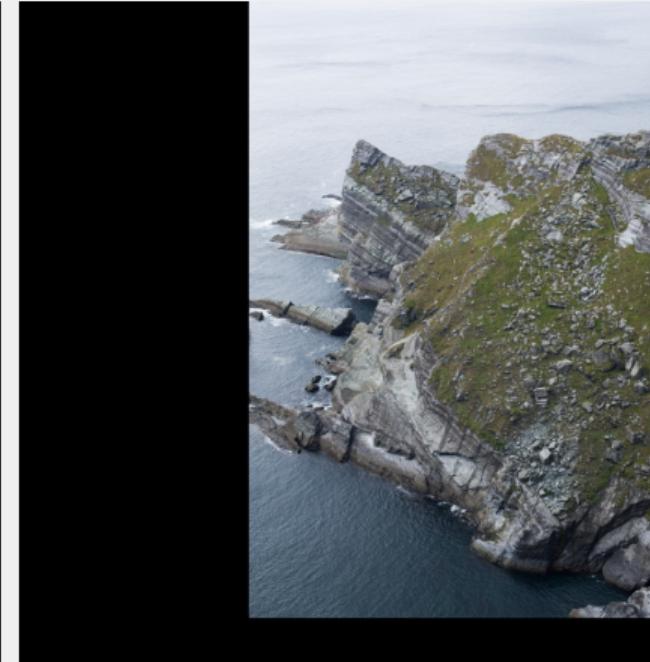
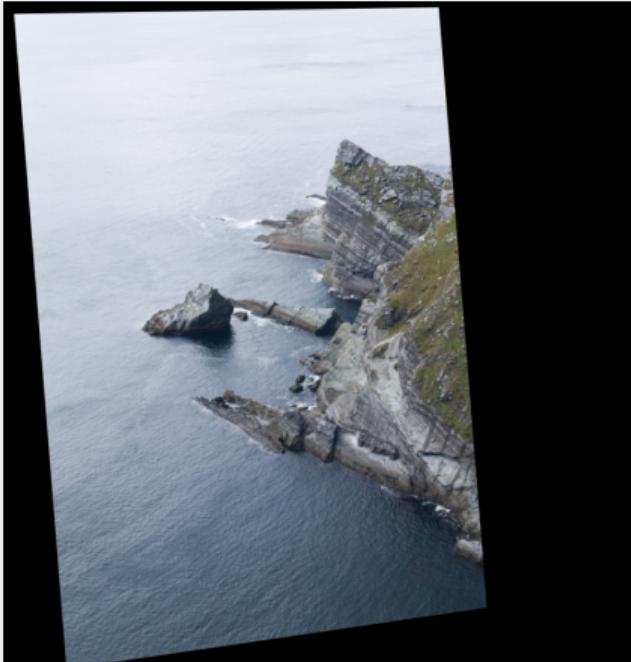
- We can use the homography to warp an image to our current field of view.
- Generate all x, y coordinates for all pixels in the reference image
- Map these to the other image using the homography
- Use bilinear interpolation to compute the value at the transformed pixel locations.
 - Code is provided in the exercise.

Warping

- You have to decide in which area to evaluate the homography
- Simple, evaluate all homographies in the same area
 - Wasteful computation
 - Simple to implement
- Warp only the valid part of each image and move the images around
 - Necessary for larger panoramas

Example – Warped to same coordinate system

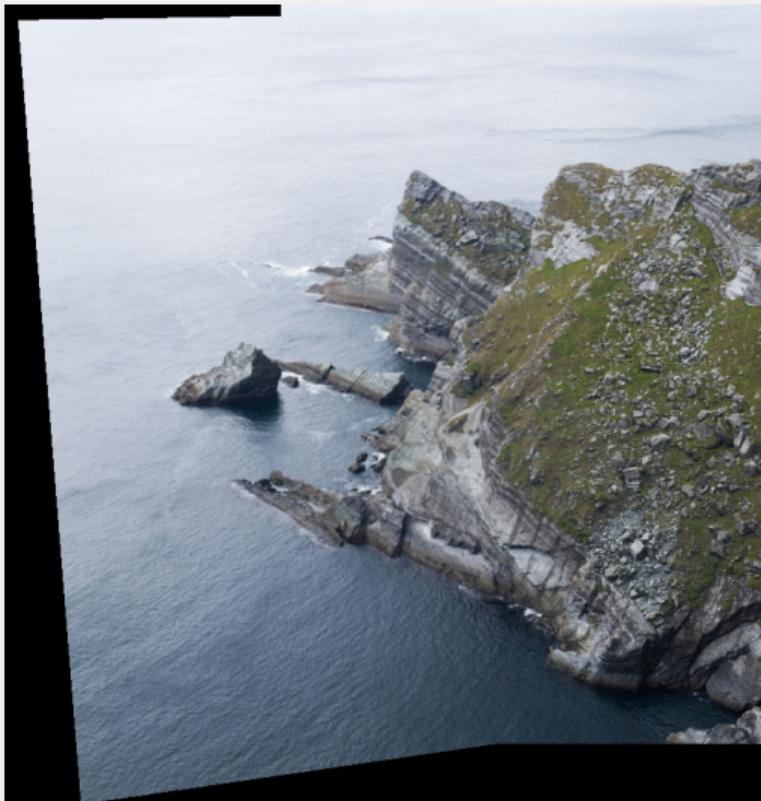
Note that x is negative.



Compositing images

- Overlap
- Average
- Median
- Graph cut

Example – Composited



Compositing – Average

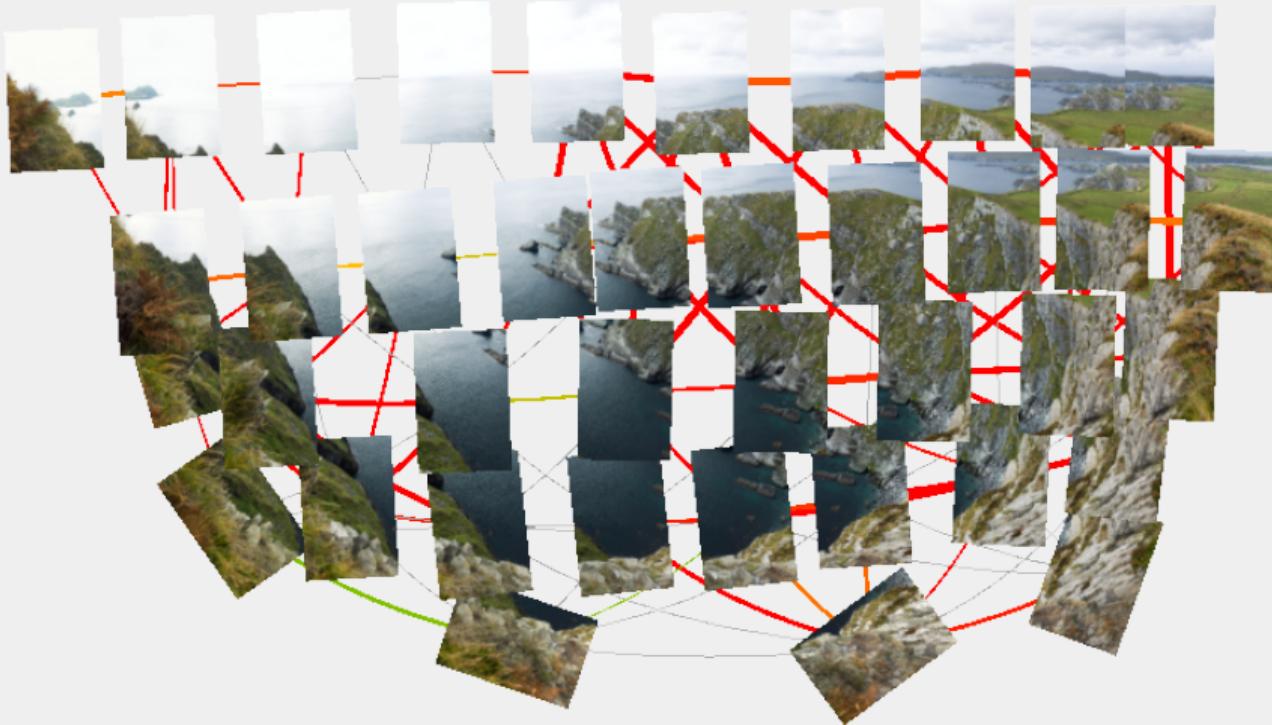


Compositing – Graph Cut



Multiple images

Multiple images



Multiple images – Final result



How do we handle more than two images?

Example with three images:

Find the homographies between each pair $H_{2 \rightarrow 1}$ and $H_{3 \rightarrow 2}$.

$$p_1 = H_{2 \rightarrow 1} p_2$$

$$p_2 = H_{3 \rightarrow 2} p_3$$

$$p_1 = \underbrace{H_{2 \rightarrow 1} H_{3 \rightarrow 2}}_{H_{3 \rightarrow 1}} p_3$$

Products of homographies are new homographies!

This principle extends to as many images as desired.

Multiple images

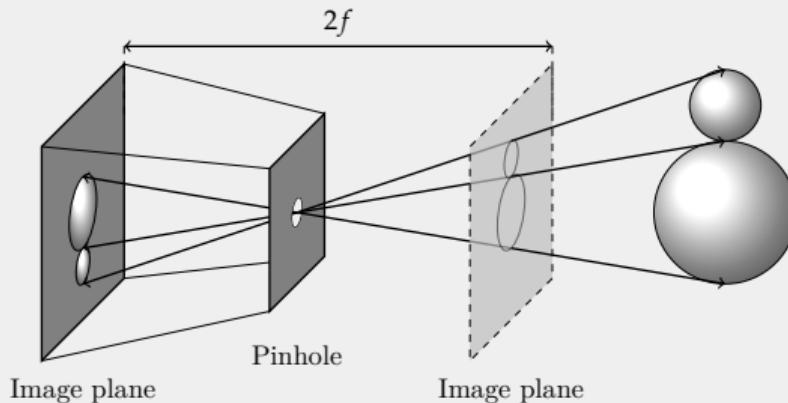
- Choose composting surface
 - I.e. which image is our baseline
 - Easy option is to choose one image to be the center
- Joint optimization

Nodal point

Capturing your own images

Which point should you actually rotate around?

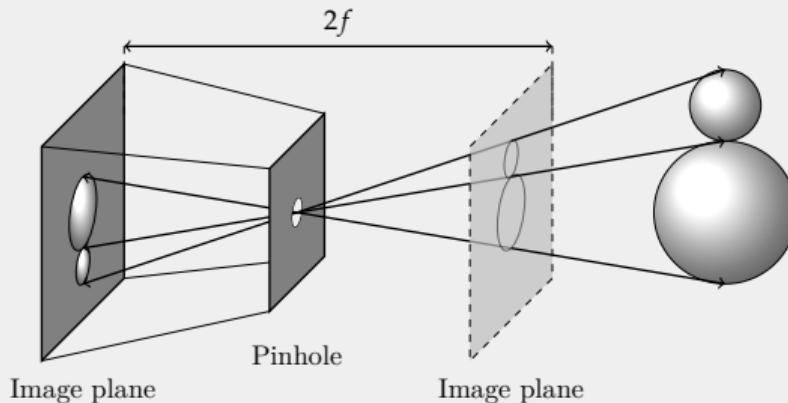
- Image sensor?
- Pinhole?



Capturing your own images

Which point should you actually rotate around?

- Image sensor?
- Pinhole?



Real world camera

Which point should you actually rotate around?



Real world camera

Which point should you actually rotate around?



It depends on the camera, usually the middle of the lens.

What happens if we rotate around the wrong point?

Real world demo[ish]

And more!

- Known focal length reduces degrees of freedom
- Other projection models
- Exposure/color correction

Learning objectives

After this lecture you should be able to:

- explain and implement RANSAC to fit homographies
- understand explain the challenges involved in stitching panoramas

Exercise time!

Visual Odometry

Morten R. Hannemose, mohan@dtu.dk

April 25, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- choose the correct decomposition of the essential matrix
- explain why the scale of the translation is unknown
- explain the Perspective- n -Point problem
- implement a simple visual odometry algorithm

Presentation topics

Decomposing the Essential Matrix

Perspective- n -Point

Putting it all together

Motivation

Let's say you're a rover on Mars...



Motivation

SLAM answers the questions:

- How is this camera moving through the world?
- What is the shape of the world?

Many applications:

- Drones
- Robotic vacuum cleaners
- Virtual reality headsets
- Augmented reality
- Autonomous cars

Many similar and related concepts

- Visual Odometry
- SLAM (Simultaneous Localization and Mapping)
- SfM (Structure from motion)

They all deal with some form of estimating a **3D map** of the world and **camera poses**, but have emphasis on different parts.

Multiple “cameras”



The unknown scale of t – Mathematical argument

$$\mathbf{E} = \mathbf{R}[t]_{\times}$$

$$0 = \mathbf{E}(st) = \mathbf{R}[t]_{\times}(st)$$

We can see that t lies in the null space of \mathbf{E} but also that it can be arbitrarily scaled

The unknown scale of t – Conceptual reason



Decomposing the Essential Matrix

Essential matrix

You have matched features between two cameras, and want to make it robust.

Estimate F or E with RANSAC.

Estimating E

- How many points are required?
- Ask yourself:
 - How many degrees of freedom does it have?
 - How many degrees of freedom does a single point fix?

Estimating E

- How many points are required?
- Ask yourself:
 - How many degrees of freedom does it have?
 - How many degrees of freedom does a single point fix?
- Five.
- Not possible with linear algorithm from five points.
- Typically estimated using Nister's five-point algorithm
 - Involves solving tenth degree polynomial
 - Is implemented in OpenCV.

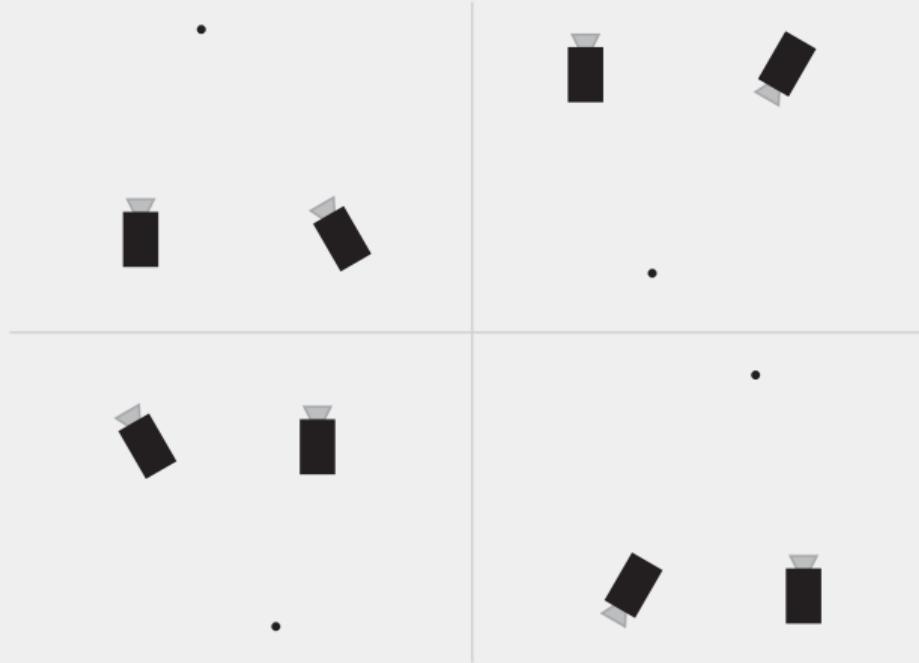
Decomposing the Essential Matrix

- The essential matrix can be computed from R and t .
 - Can we recover them from E ?

Decomposing the Essential Matrix

- The essential matrix can be computed from R and t .
 - Can we recover them from E ?
- Decomposing the essential matrix is ill posed
 - Two possible rotations
 - The sign of the translation is unknown.
- A total of four possible poses for the second camera
 - Check all four combinations
 - Choose the one with the most points in front of both cameras.

Decomposing the Essential Matrix



Back to visual odometry

We can find the pose of two cameras relative to each other

- How can we estimate the pose of a third camera?
- Using the essential matrix again will give us a new arbitrarily scaled translation

Back to visual odometry

We can find the pose of two cameras relative to each other

- How can we estimate the pose of a third camera?
- Using the essential matrix again will give us a new arbitrarily scaled translation
- Idea: Use the translation between the first two cameras to fix the scale.
- Triangulate points using the first two cameras
 - Use these 3D points to find the pose of the third camera

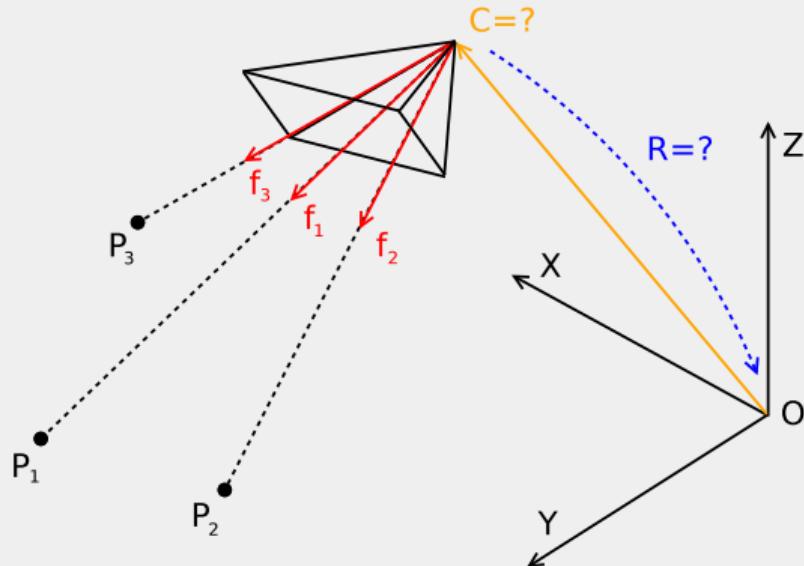
Short break

Perspective- n -Point

Perspective- n -Point (PnP)

The Perspective- n -Point (PnP) problem.

Estimating the pose of a calibrated camera from n corresponding 3D-2D correspondences.



Kneip, Laurent, Davide Scaramuzza, and Roland Siegwart. "A novel parametrization of the perspective-three-point problem for a direct computation of absolute camera position and orientation." CVPR 2011. IEEE, 2011.

PnP vs camera resectioning

In week 4 you did this for an uncalibrated camera (pest).

For an uncalibrated camera it is also called camera resectioning.

Naïve solution

- Estimate the projection matrix (\tilde{P})
- Compute $K^{-1}\tilde{P}$
- $K^{-1}\tilde{P} \approx [R \ t]$
 - R is likely not a proper rotation matrix
 - Requires many points

Perspective- n -Point (PnP)

- How many points are required?
- Ask yourself:
 - How many degrees of freedom does it have?
 - How many degrees of freedom does a single point fix?

Perspective- n -Point (PnP)

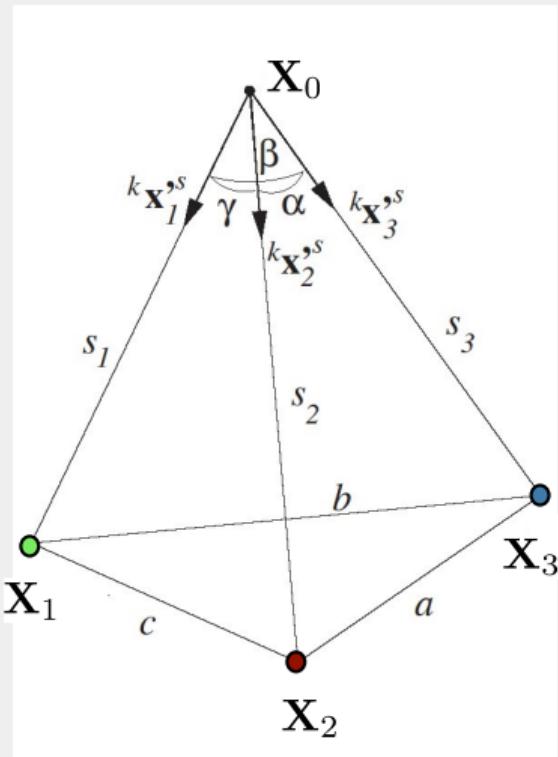
- How many points are required?
- Ask yourself:
 - How many degrees of freedom does it have?
 - How many degrees of freedom does a single point fix?
- Three correspondences are required
- This minimal case is therefore also known as P3P

P3P – Geometry

- The three 2D points give three pairwise angles
- The distances between the three 3D points give three pairwise distances

P3P – Geometry

- The three 2D points give three pairwise angles
- The distances between the three 3D points give three pairwise distances



PnP in summary

- Three correspondences generate four possible solutions, and a fourth correspondence can be used to choose the correct one
 - PnP requires 3+1 correspondences.
- Multiple algorithms exist and are implemented in OpenCV
- Add RANSAC to make it robust.

Pose vs. position

The pose of a camera is given by \mathbf{R} and \mathbf{t}

This is not the orientation and position of the camera

$$\mathbf{T}_{\text{world} \rightarrow \text{cam}} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

$$\mathbf{T}_{\text{cam} \rightarrow \text{world}} = \mathbf{T}_{\text{world} \rightarrow \text{cam}}^{-1} = \begin{bmatrix} \mathbf{R}^T & -\mathbf{R}^T \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

Pose vs. position

The orientation of a camera is given by \mathbf{R}^T

The position of a camera is given by $-\mathbf{R}^T t$

This is important in order to plot the camera

Putting it all together

Outline

1. Use the essential or fundamental matrix to estimate the pose of the second camera
2. Triangulate points using the known camera poses
3. Use PnP or camera resectioning to estimate the pose of the next camera
4. Repeat steps 2. and 3.

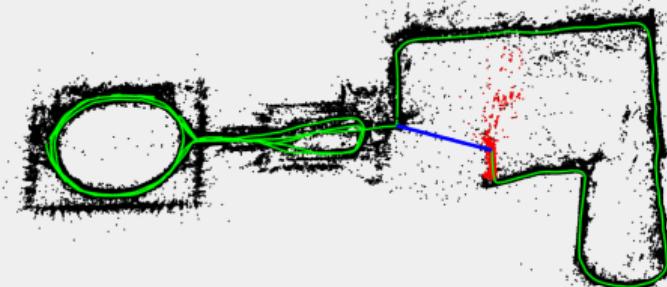
Outline

1. Use the essential or fundamental matrix to estimate the pose of the second camera
2. Triangulate points using the known camera poses
3. Use PnP or camera resectioning to estimate the pose of the next camera
4. Repeat steps 2. and 3.
5. Use (windowed) bundle adjustment

Feature tracking

- Some points can be tracked through many frames
- Keep track of them to make your model drift less

Loop closure



The exercise

- Bigger exercise (two weeks)
- Estimate essential matrix
- Estimate 3D points
- SolvePnP

Learning objectives

After this lecture you should be able to:

- choose the correct decomposition of the essential matrix
- explain why the scale of the translation is unknown
- explain the Perspective- n -Point problem
- implement a simple visual odometry algorithm

Next week

Next week: Today's exercise is bigger than usual, so you have two weeks to complete it.

You can also spend next week catching up on previous exercises.

Exercise time!

Structured light 3D scanning

Morten R. Hannemose, mohan@dtu.dk

May 9, 2025

02504 Computer vision course lectures,
DTU Compute, Kgs. Lyngby 2800, Denmark



**This lecture is being
livestreamed and recorded
(hopefully)**

Two feedback persons

Learning objectives

After this lecture you should be able to:

- explain laser line scanning
- analyse and use Gray code encoding
- analyse and use phase shift encoding

Presentation topics

Photogrammetry

Structured light

Laser line scanning

Encoding surfaces

Gray code encoding

Phase shift encoding

Notes on laser/projector calibration

Exercise: structured light

Exam information

Photogrammetry

Photogrammetry

Take a lot of pictures of a scene and use SLAM to find camera positions and 3D points. Algorithms for dense estimation of 3D points can be applied to get a full 3D scene.



Photogrammetry

Great:

- Works in daylight
- Handles textured objects

Bad:

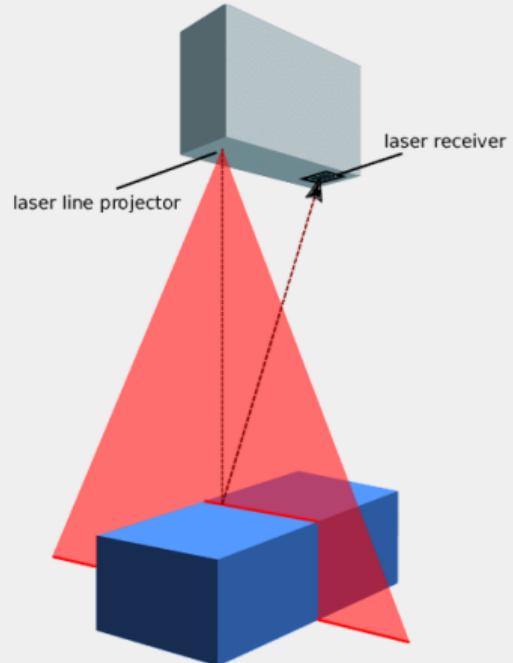
- Requires good illumination
- Requires textured objects

Structured light

Laser line scanning

The first structured light technique.

- Laser projects 3D plane of light
- Camera sees the projected line
- Laser projector and camera are calibrated



Laser line scanning

Example laser line scanner triangulation:

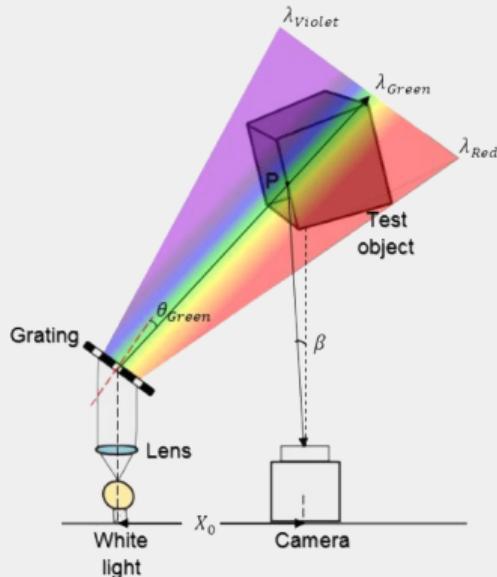
1. detect laser line in image
2. triangulate 3D point as intersection of pixel and the laser plane
3. move laser line and goto 1

Robust method, with a few drawbacks

- Requires laser calibration.
- A slow method; one triangulation line per image

Encoding surfaces

Can we encode a surface?



Yes we can!

Encoding surfaces

Possibilities:

- Continuous encoding
 - Color or intensity gradient
 - Sinusoidal (phase) shifting
- Discrete encoding
 - Binary monochrome
 - ternary RGB encoding
 - quaternary CMYK encoding
- Other encoding schemes

Encoding surfaces

- Colored encodings have problems with colored surfaces
- Intensity codings have problems with textured objects
- Continuous encoding
 - **Sinusoidal (phase) shifting**
- Discrete encoding
 - **Binary monochrome**

Continuous encoding schemes

1. For each pixel in the camera, identify the code/color
2. For each code, identify the corresponding light plane
3. Triangulate using pixel rays and the “laser plane”

Can get a 3D point for **each pixel** in the camera

Not always robust.

Discrete encoding schemes

1. For each pixel in the camera, identify the code/color
2. For each **code border**, identify the corresponding plane
3. Triangulate using pixel rays and the “laser plane”

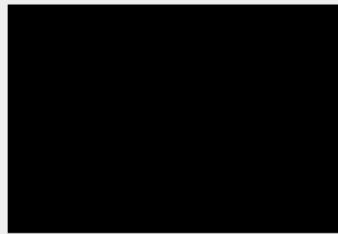
Only 3D points at **code-borders**, but usually more robust.

Binary encoding – single frame

Frame:



Inverted frame:



Frame + inverted frame makes us robust towards ambient light and varying object albedo.

Binary test: $\tau(p, p_i) = 1$ if $p > p_i$ else 0.

Single test; two regions; one border

Binary encoding – two frames

Frame 1:



Frame 2:



Inverted frame 1:

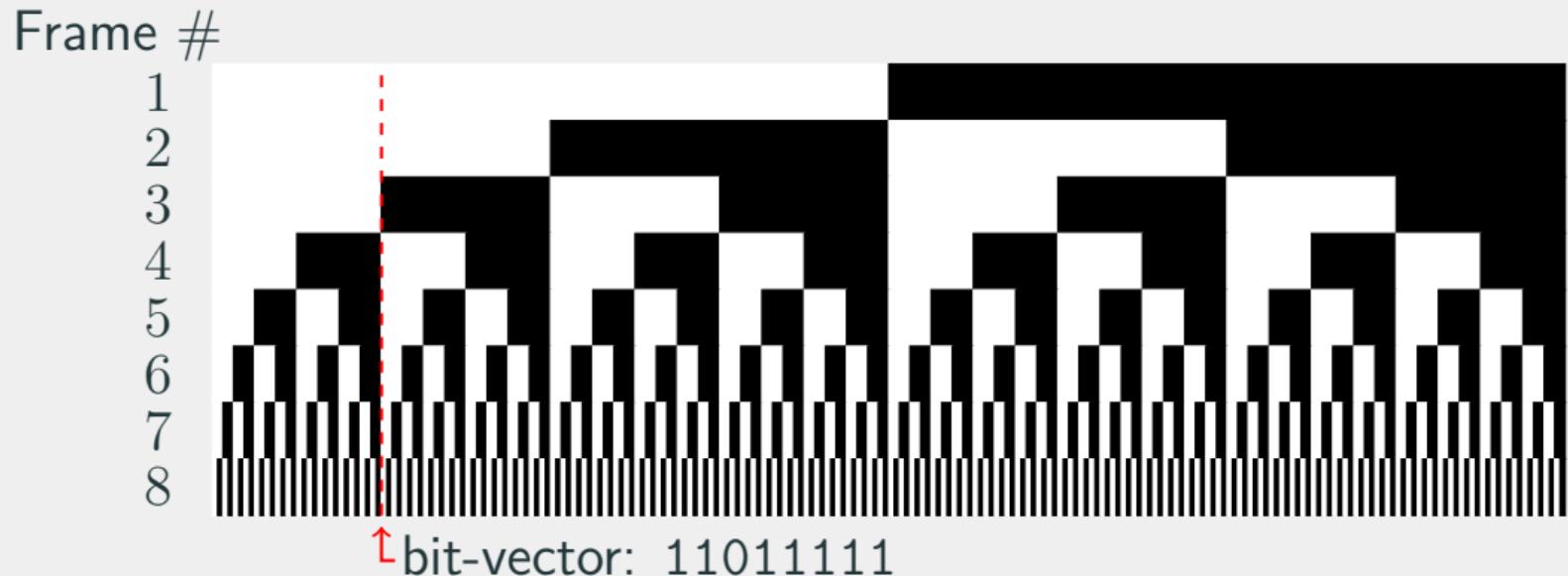


Inverted frame 2:



Two tests; four regions; three borders.

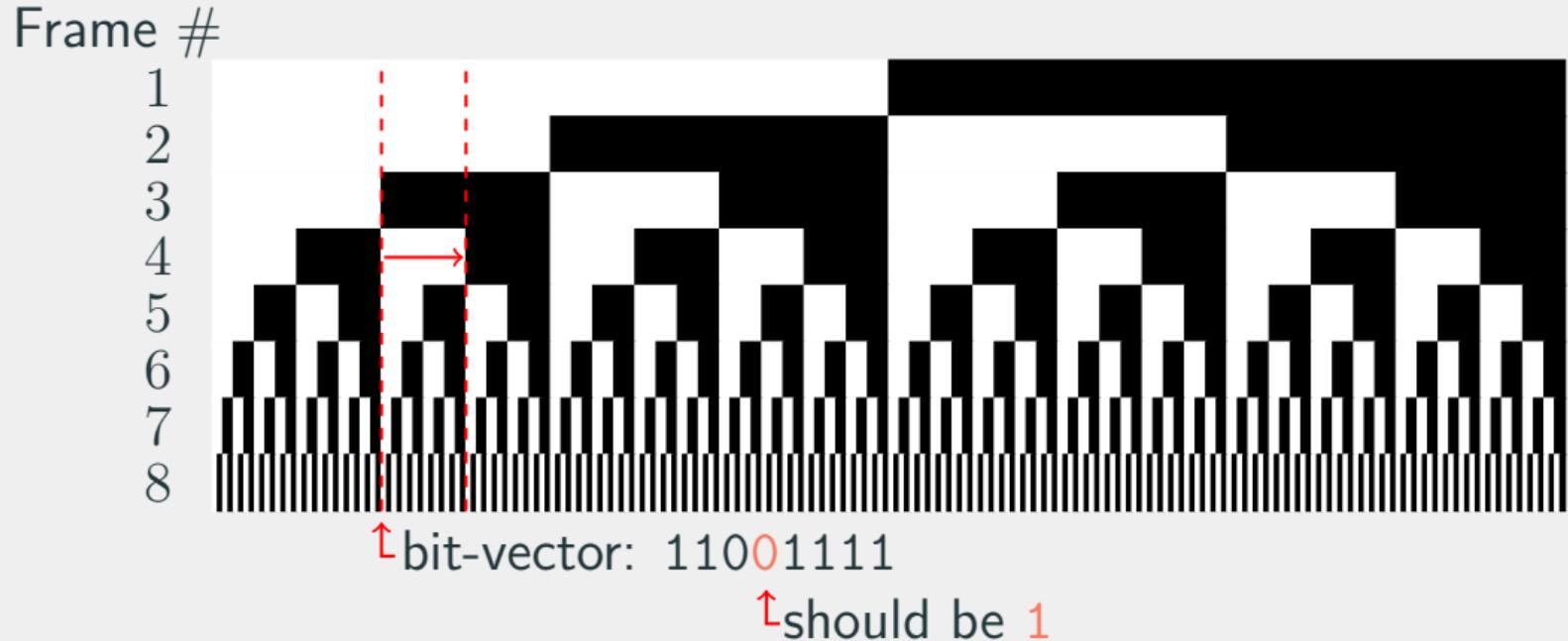
Binary encoding – multiple frames



Binary encoding

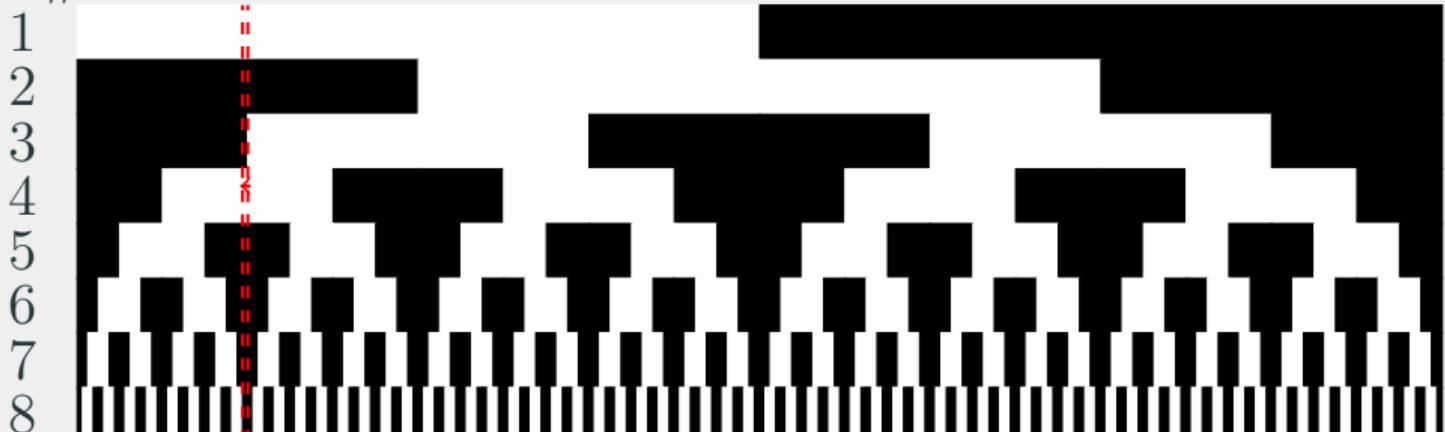
- For every frame we **subdivide the 3D volume**.
- For N frames we get 2^N unique regions.
- For N frames we get $2^N - 1$ unique borders.
- If a projector is $W = 1920$ pixels wide, we are limited to $N \leq \log_2(W) \approx 10.9$. That is a max of 10 frames in total or, 20 in total with inverted frames.

Binary encoding – border problems



Gray code encoding

Frame #



bit-vector: 10010000

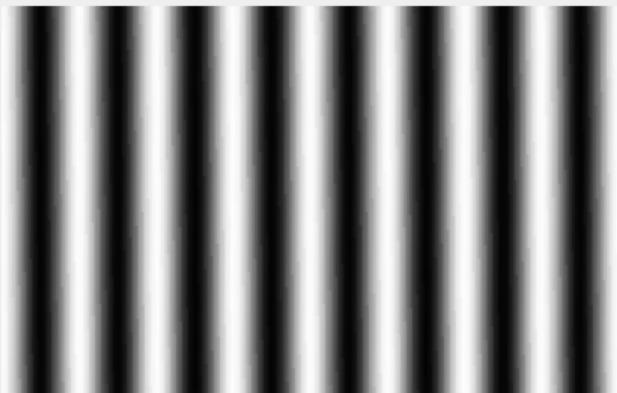
↑ should be 1

Gray codes

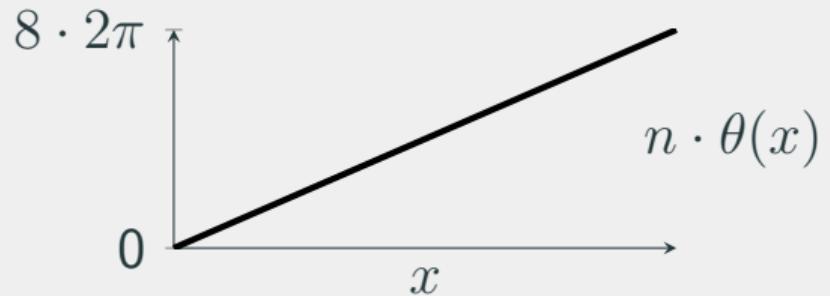
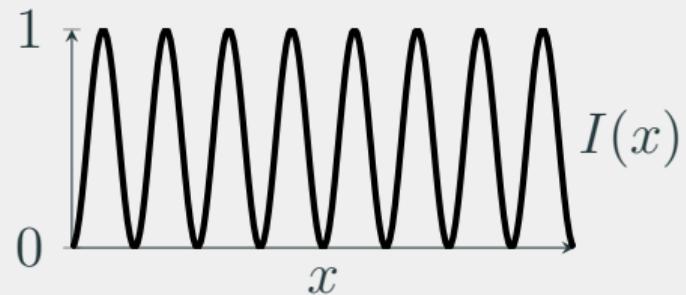
- Only one bit flip at code borders; very robust.
- Uses same number of frames as binary patterns.

Phase shift encoding

Sinusoidal waves



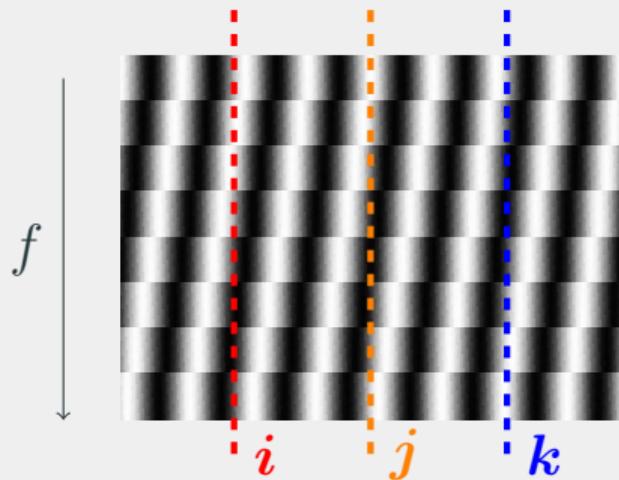
$$I(x, y) = \frac{1}{2} + \frac{1}{2} \cos(n \cdot \theta(x))$$



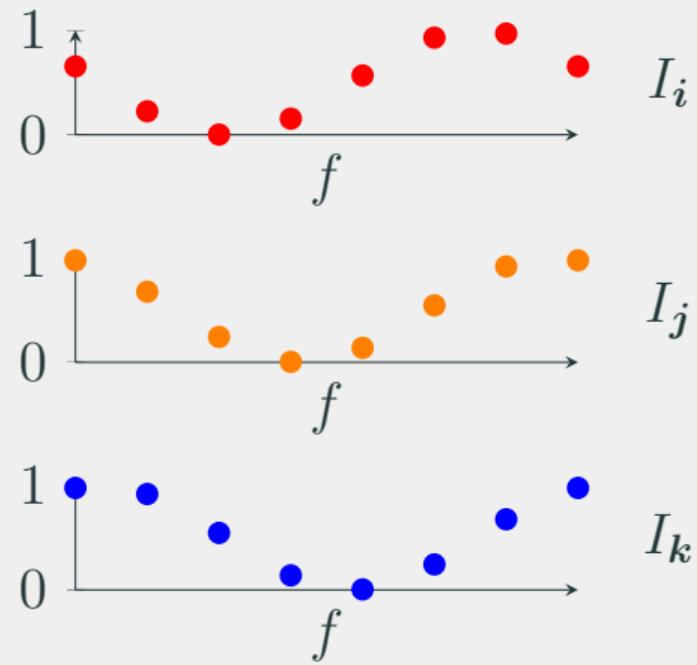
Phase as a unique code

- The pattern is monochrome; good for colored objects.
- The phase is continuous; each point in space has a unique phase-plane.
- Even with a discrete intensity projector, the pattern is approximately continuous.
- $\theta(\cdot)$ is a function of the x -coordinate of the projector
 - From now on θ is the x -coordinate of the projector (from 0 to 2π)

Phase shifting



$$I(f, \theta) = \frac{1}{2} + \frac{1}{2} \cos \left(n \cdot \theta + 2\pi \frac{f}{s} \right)$$



When projecting sinusoids and shifting them, each pixel is a sinusoid as a function of the pattern

Phase shifting method

- Phase shift exactly one wavelength in s steps.
- The phase θ corresponds to a unique projector plane.
- We can find $n \cdot \theta$ for a single pixel by fitting a sinusoid to it

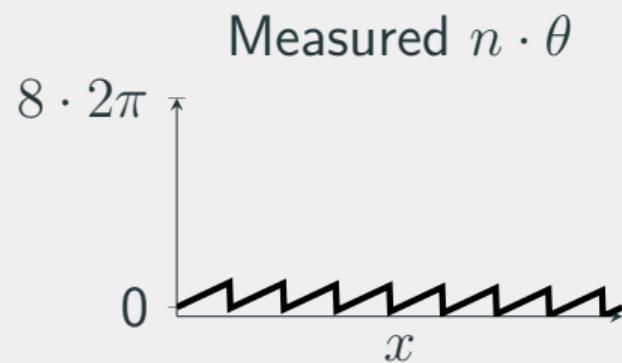
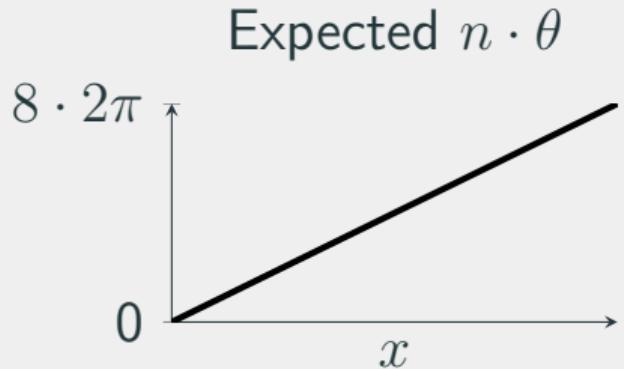
Phase shifting method

- Phase shift exactly one wavelength in s steps.
- The phase θ corresponds to a unique projector plane.
- We can find $n \cdot \theta$ for a single pixel by fitting a sinusoid to it
- This can be done with least squares (slow) or the fast Fourier transform (FFT) (fast).

$$\text{FFT}_f[I(f, \theta)] = \left\{ \frac{s}{2}, \frac{s}{2}e^{i\theta}, 0, \dots, 0 \right\}$$

The second element of the FFT is a complex number with $\theta = \text{angle}(\text{FFT}_2)$.

Phase wrapping



The phase $n \cdot \theta = \text{angle}(\text{FFT}_2)$ is wrapped to $0 \leq n \cdot \theta \leq 2\pi$.

We need to **unwrap** the recovered phase!

Heterodyne principle

Use two sinusoids:

$$\theta_1 = n_1 \cdot \theta \mod 2\pi \quad (\text{primary pattern})$$

$$\theta_2 = n_2 \cdot \theta \mod 2\pi \quad (\text{secondary pattern})$$

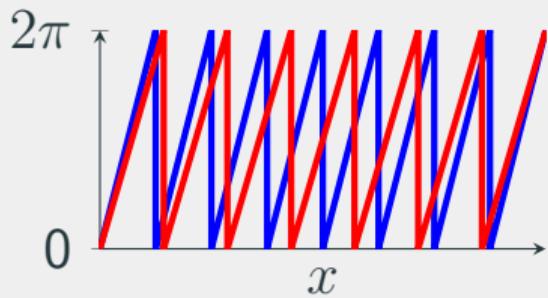
If we choose $n_2 = n_1 + 1$ we can recover θ by subtracting the primary phase from the secondary phase

$$\theta_2 - \theta_1 = n_2 \cdot \theta - n_1 \cdot \theta = \theta \mod 2\pi$$

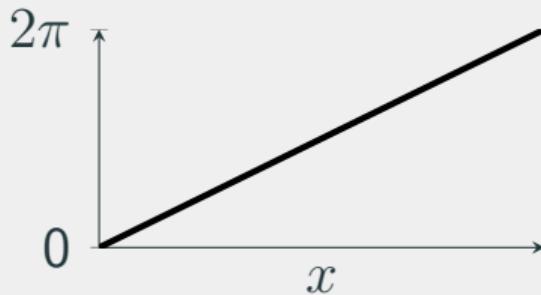
This is the heterodyne principle, and gives us the **phase cue** (θ_c).

Heterodyne principle

Measured θ_1 and θ_2



$$\theta_c = \text{mod}(\theta_2(x) - \theta_1(x), 2\pi)$$



Looks great! 😊

Questions? Short break

Unwrapping

- If the measurements of θ_1 and θ_2 are noise free, the phase cue θ_c is exactly equal to θ .
- All measurements have noise, but in this case we can improve the handling of noise substantially.

Unwrapping

To make the system more robust to noise in the measurements we can compute θ using the phase cue and the primary phase θ_1 .

The **order** counts how many times θ_1 has wrapped around

$$o_1 = \left\lfloor \frac{n_1 \cdot \theta_c - \theta_1}{2\pi} \right\rfloor$$

The rounding $\lfloor \cdot \rfloor$ makes it robust to noise. We can now estimate θ

$$\theta_{\text{est}} = \frac{2\pi o_1 + \theta_1}{n_1} \mod 2\pi$$

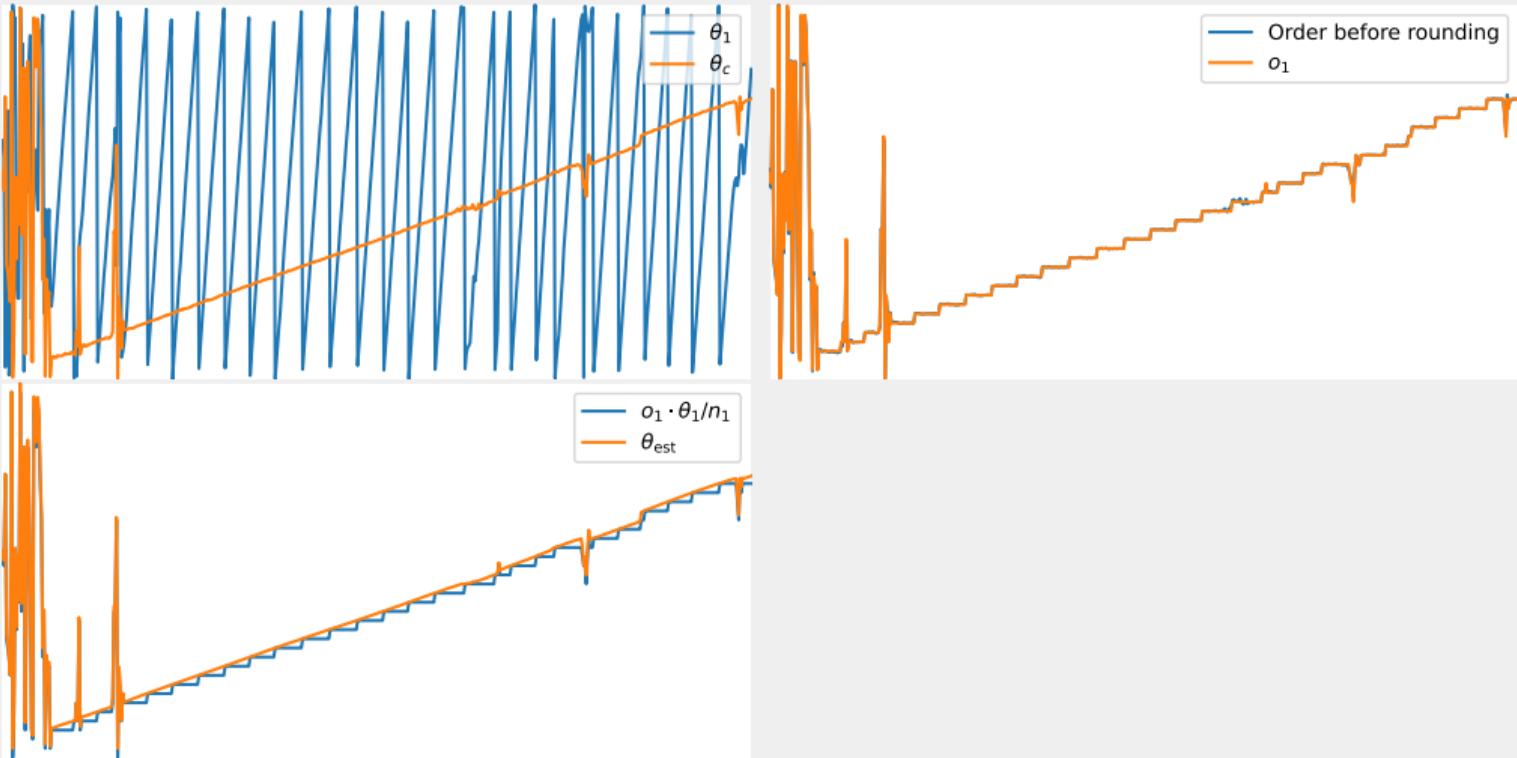
Unwrapping and phase cue – Motivation

- Why not just use the phase cue?
- $\theta_1 = \theta_1^{\text{true}} + \epsilon_1$
- $\theta_2 = \theta_2^{\text{true}} + \epsilon_2$
- $\theta_c = \theta^{\text{true}} - \epsilon_1 + \epsilon_2$

Unwrapping and phase cue – Motivation

- Why not just use the phase cue?
- $\theta_1 = \theta_1^{\text{true}} + \epsilon_1$
- $\theta_2 = \theta_2^{\text{true}} + \epsilon_2$
- $\theta_c = \theta^{\text{true}} - \epsilon_1 + \epsilon_2$
- The error of the phase cue is linear in the errors of θ_1 and θ_2
- θ has an error of $\frac{\epsilon_1}{n_1}$, which is much lower

Unwrapping – examples from today's exercise.



Examples from row 400 of camera 0.

Phase shift encoding

1. Project two sets of sinusoidals.
2. For each pixel, find the wrapped phases θ_1 and θ_2 .
3. Calculate the phase cue $\theta_c = \theta_2 - \theta_1 \bmod 2\pi$.
4. Calculate the primary order $o_1 = \left\lfloor \frac{n_1 \cdot \theta_c - \theta_1}{2\pi} \right\rfloor$
5. Calculate the theta $\theta_{\text{est}} = \frac{2\pi o_1 + \theta_1}{n_1}$
6. For each phase, find the corresponding projector plane.
7. Triangulate pixel rays and the projector plane.

Phase shift encoding

- Very robust since the FFT is also a low-pass filter.
- Potentially one triangulation per camera pixel.
- Easy to make more precise; more projected frames gives a higher precision.

Notes on laser/projector calibration

Laser/projector calibration

- Is needed for many structured light methods.
- Requires correspondences between the laser/projector and the real world.
- Often we use a calibrated camera to relate projected lines/pixels to world coordinates.
- Not ideal for projectors, as they usually have poor quality lenses.

However, do we need to calibrate the laser/projector?

“Passive” structured light

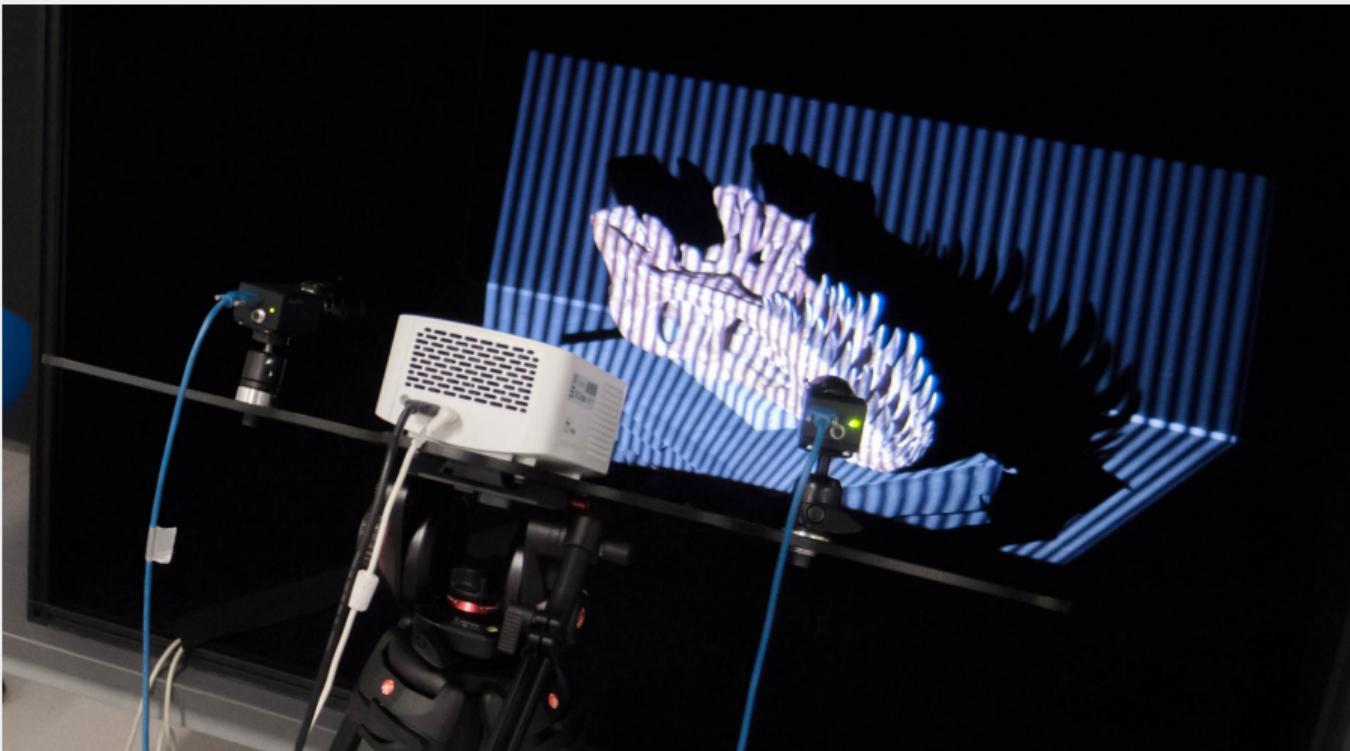
- Use a stereo camera set up and a laser/projector.
- The laser/projector encodes the object surface, but is not calibrated.
 - Accurately calibrating a projector is harder than calibrating a camera
- Instead of laser/projector plane triangulation, we use epipolar lines in the cameras.
- Requires that the projected planes are not parallel with the epipolar planes.

Exercise: structured light

Structured light exercise

- Scan using phase shifting
- You will perform:
 - Image rectification
 - Phase decoding (and unwrapping)
 - Masking
 - Matching
 - Triangulation

Casper the baby T-Rex

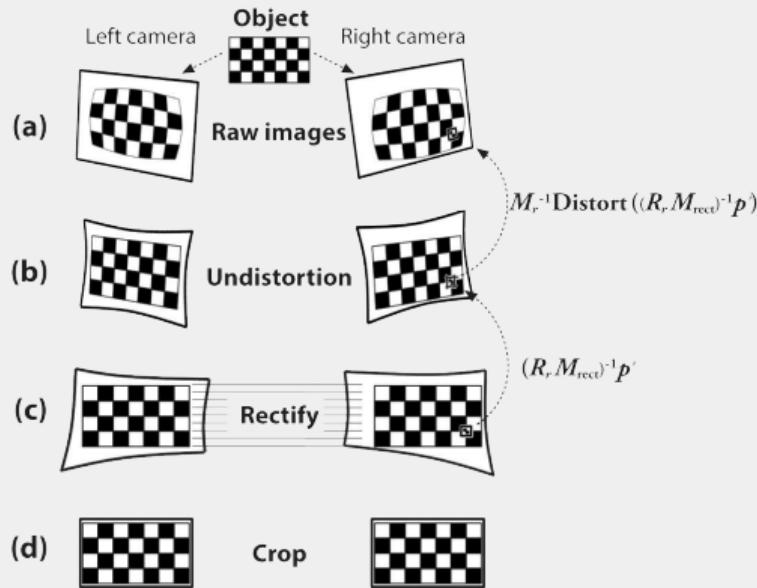


Casper the baby T-Rex



Rectified cameras

Cameras are (virtually) made parallel to each other.



Masking

We need to only match phases, which is efficient and robust.

With rectified images, epipolar lines are the corresponding rows.

Live demo[ish]

Further courses

- Are you interested in working more with images?
- 02506 Advanced Image Analysis (Spring)
- 02516 Introduction to Deep Learning in Computer Vision (Fall)
- 02501 Advanced Deep Learning in Computer Vision (Spring)

Learning objectives

After this lecture you should be able to:

- explain laser line scanning
- analyse and use Gray code encoding
- analyse and use phase shift encoding

Exam information

Exam information

Exam time and place: www.eksamensplan.dtu.dk

- Multiple choice exam
 - No negative points for wrong answers
- Questions are a mix of understanding and calculation in Python.
 - You can use your computer and all notes.
 - No internet access during exam.
 - No use of LLMs/generative AI.
 - You can [download](#) the documentation for OpenCV to your computer.

Quiz awards

Exercise time!