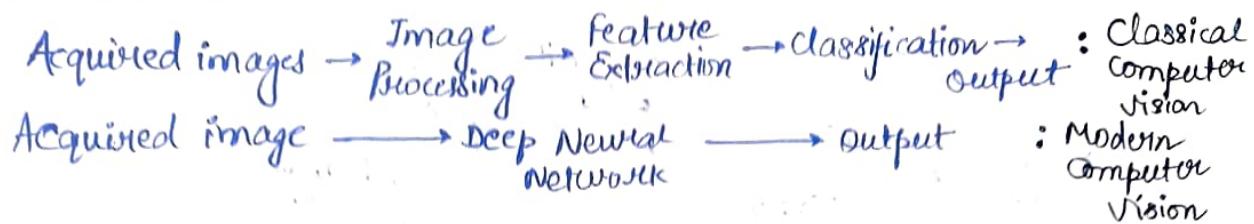


# COMPUTER VISION



## Image Formation

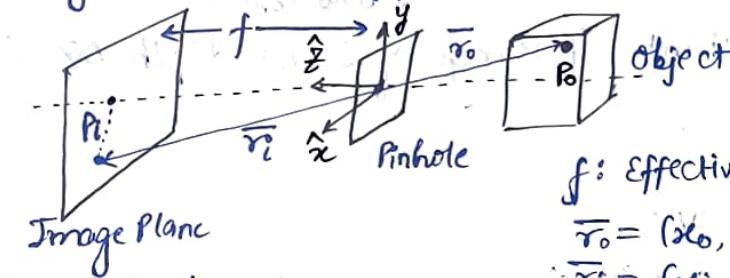
### Geometry of Image Formation

↳ Mapping b/w image and world coordinates.

- ① Pinhole camera Model
- ② Projective geometry
- ③ Projective Matrix.

### Image Formation in camera:

↳ Image: Projection of 3D scenes onto 2D plane.



$f$ : effective focal length

$$\vec{r}_o = (x_o, y_o, z_o)$$

$$\vec{r}_i = (x_i, y_i, f)$$

Using similar  $\Delta$ s,

$$\frac{\vec{r}_i}{f} = \frac{\vec{r}_o}{z_o} \Rightarrow \frac{x_i}{f} = \frac{x_o}{z_o}, \frac{y_i}{f} = \frac{y_o}{z_o}$$

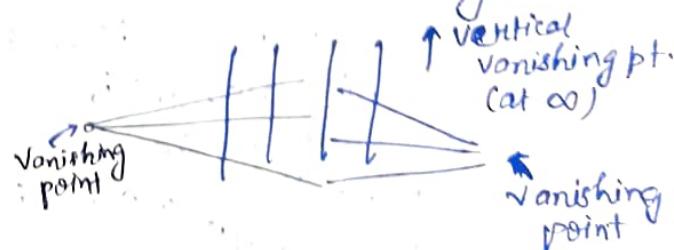
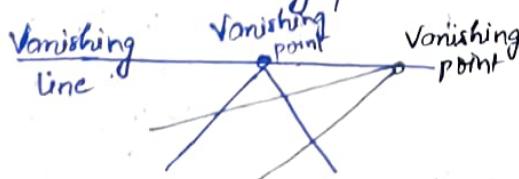
Pinhole: barrier  
→ Reduce blurring  
→ Opening known aperture.

↳ Lost: Length, angles

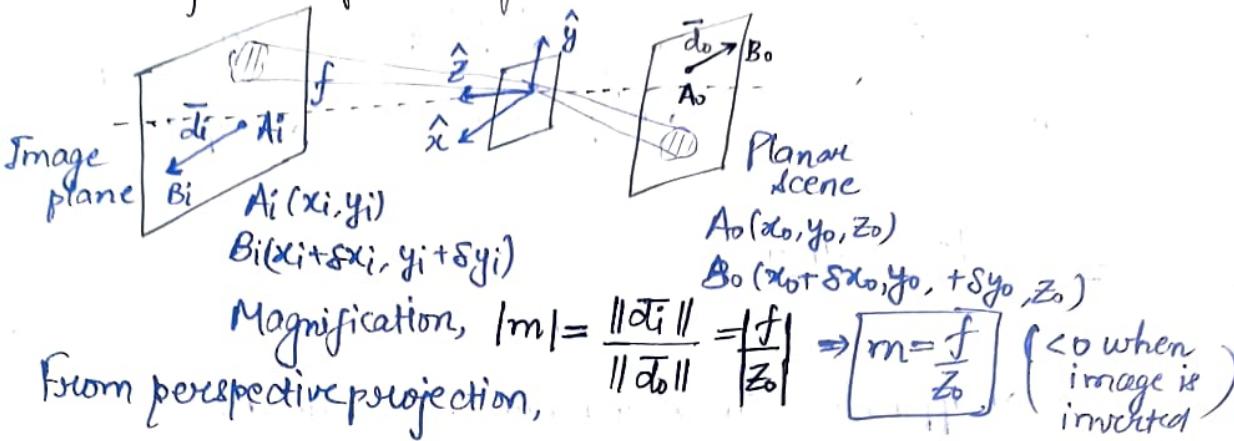
Preserved: Straight lines are still straight.

## Vanishing Points and Lines

↳ Parallel lines in the world intersect in the image at a 'vanishing point'.



## Perspective projection of a straight line



$$\frac{x_i}{f} = \frac{x_0}{z_0}, \quad \frac{y_i}{f} = \frac{y_0}{z_0}$$

$$\Rightarrow \frac{\delta x_i}{f} = \frac{\delta x_0}{z_0}, \quad \frac{\delta y_i}{f} = \frac{\delta y_0}{z_0}$$

$$\therefore |m| = \frac{\|\vec{d}_i\|}{\|\vec{d}_0\|} = \frac{\sqrt{\delta x_i^2 + \delta y_i^2}}{\sqrt{\delta x_0^2 + \delta y_0^2}} = \frac{|f|}{|z_0|}$$

→  $m \approx$  constant, if range of scene depth  $\Delta Z \ll$  average scene depth  $\bar{Z}$ .

$$\frac{\text{Area}_i}{\text{Area}_0} = m^2$$

### Challenges in CV:

- Perspective illusion
- Viewpoint variation
- Illumination
- Motion
- Perception vs. Measurement
- Intradclass variability

### Applications:

- Optical character recognition (OCR)
- Face detection, smile detection
- Vision-based biometrics, facial login
- Object recognition
- 3D from images
- Special effects: Motion capture
- Interactive games: Kinect
- Sports
- Medical imaging
- Autocars
- Industrial robots
- Vision in space
- Mobile robots
- Augmented reality & Virtual reality.

## Pinhole Camera and Perspective Projection

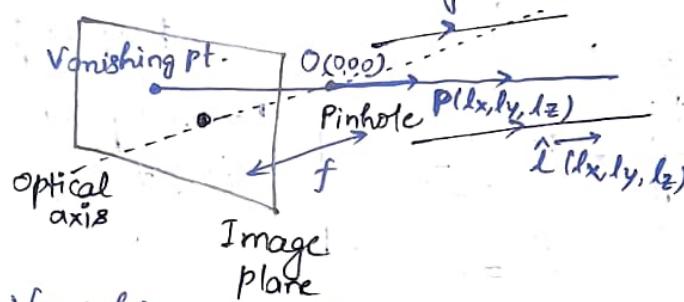
Perspective projection: Mapping from 3D onto 2D.

Pinhole camera: Imaging device which captures accurately the geometry of perspective projection.

Camera & world coordinate system

↳ In general, not aligned.

Vanishing Point: location depends on the orientation of parallel straight lines.



Vanishing pt. of line: Projection of point P

$$(x_{vp}, y_{vp}) = \left( f \frac{x}{l_z}, f \frac{y}{l_z} \right).$$

Pinhole → tiny

↳ cause diffraction, if it's too tiny.

Ideal pinhole diameter:  $(d \approx 2\sqrt{f\lambda})$

f: effective focal length  
 $\lambda$ : wavelength

## Lenses

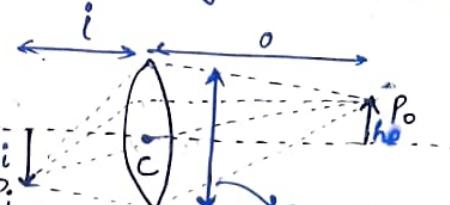
↳ Same projection as pinhole, but gather more light.  
Gaussian lens (thin lens) law:

$$\frac{1}{i} + \frac{1}{o} = \frac{1}{f}$$

f: focal length

i: image/object distance

If  $o = \infty$ ,  $f = i$ .



D: aperture/diameter of lens

Reduce/increase to control image brightness.

$$\text{Magnification: } m = \frac{h_i}{h_o} = \frac{i}{o}$$

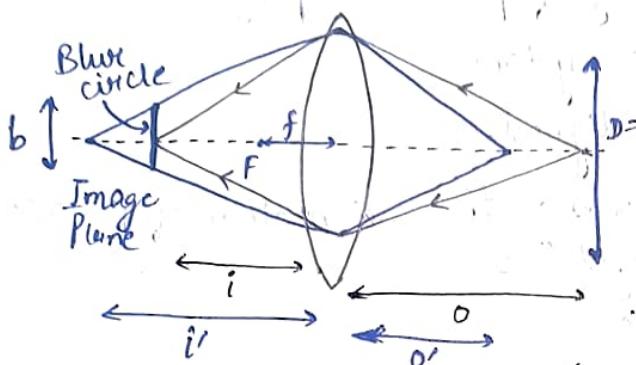
For two lens system,  $m = \frac{i_2}{o_2} \cdot \frac{i_1}{o_1}$

$$\text{Aperture, } D = f/N$$

$$f\text{-Number, } N = f/D$$

- lens distortion (vignetting)
- chromatic aberration
- geometric distortion
  - ↳ Radial, tangential

## Lens Defocus



Defocused point:

$$\frac{1}{i'} + \frac{1}{o'} = \frac{1}{f} \Rightarrow i' = \frac{o'f}{o'-f}$$

$$\Rightarrow i' - i = \frac{f}{o'-f} \cdot \frac{f}{o'-f} \cdot (o - o')$$

$$\therefore b = Df \left| \frac{(o - o')}{o'(o - f)} \right|$$

From similar  $\Delta s$ ,

$$\frac{b}{D} = \frac{|i' - i|}{|i'|}$$

Blur circle diameter,

$$b = \frac{D}{i'} |i' - i|$$

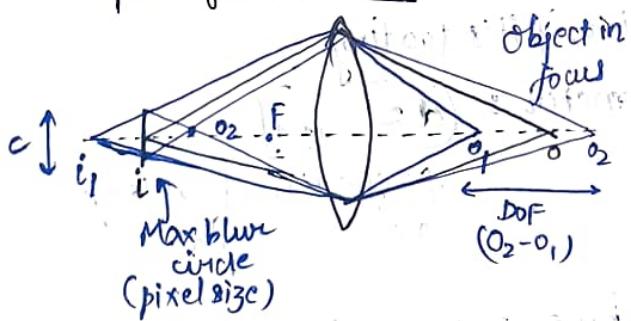
$$b \propto D \propto \frac{1}{N}$$

Focused point

$$i = \frac{of}{o-f}$$

Move image plane or lens or both for focusing.

## Depth of Field (DoF)

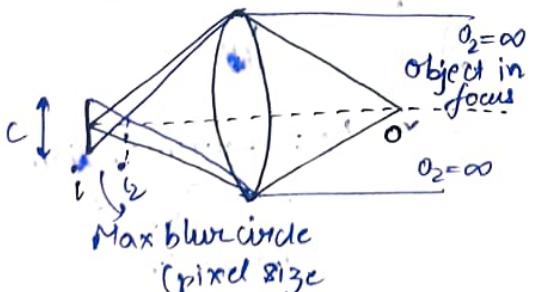


$o_1, o_2$ : nearest & farthest distances for blur circle to be max<sup>m</sup>c.

$$c = \frac{f^2(o - o_1)}{N o_1(o - f)} = \frac{f^2(o_2 - o)}{N o_2(o - f)}$$

$$\text{DoF: } \frac{o_2 - o_1}{f^2 - c^2 N^2 (o - f)^2}$$

## Hypofocal distance



The closest distance  $o = h$ , the lens must be focused to keep object at  $\infty$  ( $o_2$ ) acceptably sharp (blur circle  $\leq c$ ).

hypofocal distance

$$h = \frac{f^2}{Nc} + f$$

→ large aperture (small f no.)

↳ Bright image or short exposure time

↳ Shallow depth of field

→ small aperture (large f no.)

↳ Dark image or long exposure time

↳ Large DoF.

Allow camera rotation:

$$x = K[R + t]x \Rightarrow w \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha & s & u_0 \\ 0 & p & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \gamma_{11} & \gamma_{12} & \gamma_{13} & t_x \\ \gamma_{21} & \gamma_{22} & \gamma_{23} & t_y \\ \gamma_{31} & \gamma_{32} & \gamma_{33} & t_z \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

Degrees of freedom : (5) (6)

## Projective geometry:

- ① A circle may appear as an ellipse.
- ② Parallel lines may meet at a finite point.
- ③ A rectangle may appear as a parallelogram.  
→ Angles, distances, ratio of distances are not preserved.
- $(Kx, Ky, Kz)$  corresponds to  $\left(\frac{Kx}{K}, \frac{Ky}{K}\right) = (x, y)$   
in homo. coord. in cartesian coord.
- $(Kx, Ky, Kz)$  are equivalent for all  $K$ 's.
- Point at infinity :  $(x, y, z)$  where  $z=0$   
i.e.,  $(\frac{x}{0}, \frac{y}{0})$ .

## Homogeneous Point

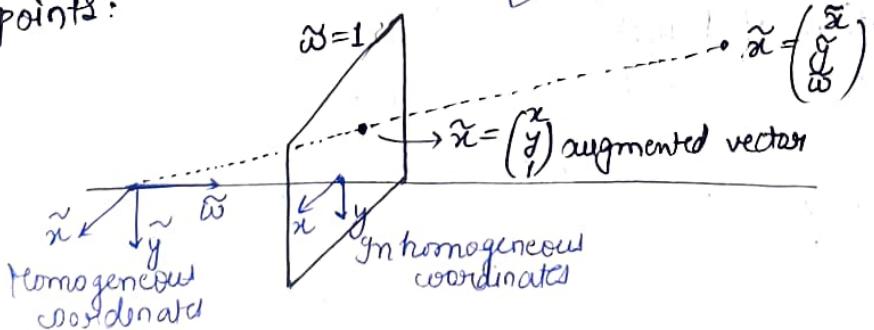
$$x = \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \in \mathbb{R}^3$$

$$\tilde{x} = \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \in \mathbb{P}^2$$

where  $\mathbb{P}^2 = \mathbb{R}^3$  is called projective space

[Except the element  $(0, 0, 0)$ ]

2D-points:



# CAMERA

→ CCD & CMOS camera: Convert light into electrical signals to produce images.

CCD Sensor: Have passive pixels that transfer charge to a single amplifier

CCD Camera: → Have better image quality, especially in low-light settings.

→ Have Charged Coupling Device

→ Photons → electrons → no. of → bits

→ Have alternative RGB pixels that capture the respective colors.

## CCD Camera Sensor

- Have passive filters that transfer the charge to a single amplifier
- Charge coupled device.
- Each charge in each pixels are transferred using horizontal and vertical shift registers, then converted to voltage and amplified.
- Not possible to integrate peripherals like timers and ADC, in the main sensor; required additional chip and thus large size.
- Requires different power supplies, typically 7–10 V, thus consuming more power.
- Low processing speed as each pixel is processed one-by-one.
- Large fill factor, noise and sensitivity.
- Better camera image quality.
- Use global shutter.
- Suffer from smear.
- Image distortion due to blooming.

## CMOS Camera Sensor

[SAURABH KUMAR | SC22B146]

- Have active pixels that amplify the charge in each pixel.
- Complementary Metal Oxide Semiconductor.
- Charge to voltage conversion and amplification are carried out in the pixels itself.
- Fabrication procedure is very similar to that of IC, it is possible to integrate peripheral components into single chip, thus, possible to have camera on chip or SoC.
- Requires single power supply, typically 3.3–5 V, thus consuming less power.
- High speed as processing is carried out in the pixel itself.
- Low fill factor, noise and sensitivity.
- Cheaper, energy efficient
- Use rolling shutter (R.S. artifacts).
- Affected by skew, wobble & partial exposure.

# Homogeneous Coordinate System

Homogeneous coordinates:

$$(x, y) \Rightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Homo. image  
world.

$$(x, y, z) \Rightarrow \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

Homo. scene  
coord.

[Converting to  
homogeneous  
coordinates]

$$\begin{bmatrix} x \\ y \\ w \end{bmatrix} \Rightarrow (x/w, y/w)$$

$$\begin{bmatrix} x \\ y \\ z \\ w \end{bmatrix} \Rightarrow (x/w, y/w, z/w)$$

[Converting from  
homogeneous  
coordinates]

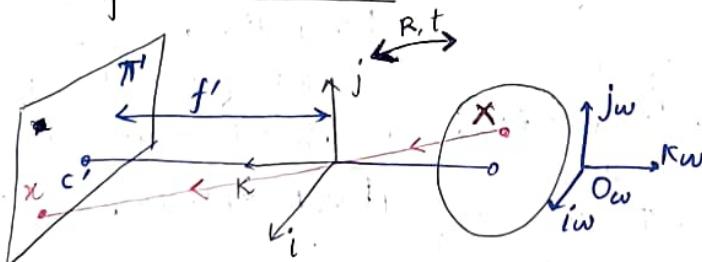
Projective space  $\rightarrow \mathbb{P}^2$

$$\begin{pmatrix} kx \\ ky \\ k \end{pmatrix} \Rightarrow \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \quad [\text{Invariant to scaling}]$$

$$\left( \frac{kx}{k}, \frac{ky}{k} \right) \rightarrow (x, y) \in \mathbb{R}^2, \quad [\text{Cartesian}]$$

$\rightarrow$  Point in Cartesian is ray in Homogeneous.

Camera Projection Matrix:



$$x = K[R \ t]X$$

Intrinsic assumption:

↳ Unit skew ratio

↳ Optical center at (0,0)

↳ No skew

Extrinsic assumptions:

↳ No rotation

↳ camera at (0,0,0)

x : image coordinates: (u, v, 1)

K: Intrinsic Matrix (3x3)

R: Rotation (3x3)

t: Translation (3x1)

X: World coordinates (x, y, z, 1)

$$x = K[I \ 0]X \Rightarrow {}^w\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

## 2-D Line:

(21-01-2025)

Parametric form of line:

$$ax + by + c = 0$$

$$\begin{bmatrix} a & b & c \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = 0 \Rightarrow l^T x = 0 ; l = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

$$(x_1, x_2, x_3) \Rightarrow \left( \frac{x_1}{x_3}, \frac{x_2}{x_3} \right)$$

Homogeneous space      Non-homo.

For Two lines  $l$  and  $l'$  intersect at  $x$ .

$$\Rightarrow \boxed{l \times l'} \quad (\text{Cross Product})$$

$$\Rightarrow l^T x = 0$$

$$\Rightarrow l^T(l \times l') = 0$$

Eg

$$y=1 \rightarrow l' = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

$$l = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}$$

$$l \times l' = \begin{vmatrix} i & j & k \\ 1 & 0 & 1 \\ -1 & 0 & 1 \end{vmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \equiv (1, 1)$$

- $ax+by+c=0$  is same as  $(ka)x+(kb)y+(kc)=0$   $\forall k \neq 0$
- $(0, 0, 0)^T \rightarrow$  not a line
- $(a, b, c)^T, k(a, b, c)^T \rightarrow$  same line
- Point  $x=(x, y)^T$  lies on line  $l = (a, b, c)^T$  iff  $ax+by+c=0$   
i.e.,  $(x, y, 1)(a, b, c)^T = 0$

→ Point at infinity: Ideal point

$$\hookrightarrow (x, y, 0) = (x_0, y_0)$$

Parallel Lines

Eg.  $x=1$  ( $l$ )  
 $x=2$  ( $l'$ )

$$l = [1 \ 0 \ -1]^T$$

$$l' = [1 \ 0 \ -2]^T$$

$$l \times l' = \begin{vmatrix} i & j & k \\ 1 & 0 & -1 \\ 1 & 0 & -2 \end{vmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

Line at infinity:  $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = l_\infty$

$$\boxed{l^T \cdot l_\infty = 0} \Rightarrow (0 \ 0 \ 1) \begin{pmatrix} x \\ y \\ 0 \end{pmatrix} = 0$$

Ideal point

Line joining two points  $x_1$  and  $x_2$ :  $\ell = x_1 \times x_2$

Cross Product:

Convert one vector to skew:  $[x]_x = \begin{bmatrix} 0 & -3 & y \\ 3 & 0 & -x \\ -y & x & 0 \end{bmatrix}$

 $\ell_1^T = (0 \ 1 \ -1) ; \ell_2^T = (1 \ 0 \ -2)$

$$\ell_1 \times \ell_2 = [\ell_1]_x \ell = \begin{pmatrix} 0 & 1 & 1 \\ -1 & 0 & 0 \\ -1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix}$$

Parallel lines:

$$(1 \ 0 \ -1) \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$

$$(1 \ 0 \ -2) \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$

$$\ell_1 \times \ell_2 = [\ell_1]_x \ell_2 = \begin{pmatrix} 1 & 0 & 1 \\ -1 & 0 & -1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

$$(\overbrace{\ell_2}^{(0 \ 0 \ 1)^T} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}) = 0$$

Conics

$$\text{Circle: } x^2 + y^2 - 1 = 0 \quad (x/3)^2 + (y/3)^2 - 1 = 0 \quad x^2 + y^2 - 3^2 = 0$$

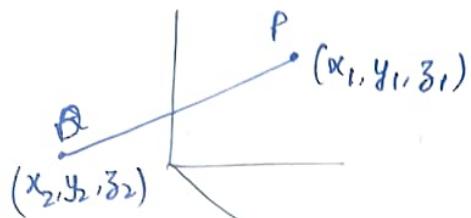
$$\text{Parabola: } x^2 - y = 0 \quad \rightarrow (x/3)^2 - (y/3) = 0 \quad \overrightarrow{xz}^2 - yz = 0$$

$$\text{Hyperbola: } x^2 - y^2 - 1 = 0 \quad (x/3)^2 - (y/3)^2 - 1 = 0 \quad x^2 - y^2 - 3^2 = 0$$

In general:

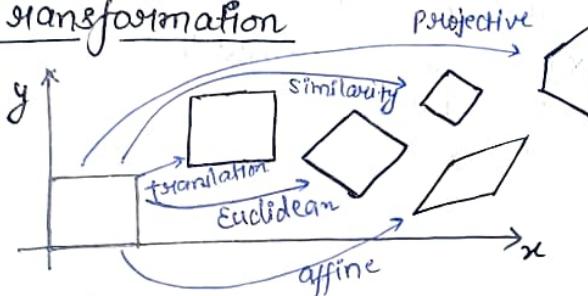
$$ax^2 + by^2 + cxy + dxz + eyz + fz^2 = 0$$

## 3D Line



Parametric form:  $(1-\lambda)P + \lambda Q$

## 2D Transformation



Translation: 2 DoF

$$x' = x + t \Leftrightarrow \bar{x}' = \begin{pmatrix} I & t \\ 0^T & 1 \end{pmatrix} \bar{x}$$

Euclidean: 2D Translation + 2D Rotation ; 3 DoF

$$x' = Rx + t \Leftrightarrow \bar{x}' = \begin{pmatrix} R & t \\ 0^T & 1 \end{pmatrix} \bar{x}$$

Similarity: 2D Translation + Scaled 2D Rotation ; 4 DoF

$$x' = sRx + t \Leftrightarrow \bar{x}' = \begin{bmatrix} sR & t \\ 0^T & 1 \end{bmatrix} \bar{x}$$

Affine: 2D Linear Transformation ; 6 DoF

$$x' = Ax + t \Leftrightarrow \bar{x}' = \begin{bmatrix} A & t \\ 0^T & 1 \end{bmatrix} \bar{x}$$

Projective: Homography ; 8 DoF

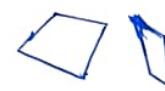
$$\tilde{x}' = \tilde{H}\tilde{x} \quad (\bar{x} = \frac{1}{\tilde{\omega}}\tilde{x})$$

## 2D Transformation Hierarchy:



Projective  
(8 DoF)

$$\begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix}$$



Affine  
(6 DoF)

$$\begin{bmatrix} a_{11} & a_{12} & tx \\ a_{21} & a_{22} & ty \\ 0 & 0 & 1 \end{bmatrix}$$



Similarity  
(4 DoF)

$$\begin{bmatrix} s\alpha_{11} & s\alpha_{12} & tx \\ s\alpha_{21} & s\alpha_{22} & ty \\ 0 & 0 & 1 \end{bmatrix}$$



Euclidean  
(3 DoF)

$$\begin{bmatrix} r_{11} & r_{12} & tx \\ r_{21} & r_{22} & ty \\ 0 & 0 & 1 \end{bmatrix}$$



## (3D) 2D Transformation:

### Transformation

Translation

| <u>Transformation</u> | <u>Matrix</u>                         | <u># DoF</u> | <u>Preserves</u> | <u>Icon</u> |
|-----------------------|---------------------------------------|--------------|------------------|-------------|
| Translation           | $[I \ t]_{2 \times 3}^{(3 \times 4)}$ | 2 (3)        | Orientation      |             |

Rigid  
(Euclidean)

|                      |                                       |       |        |  |
|----------------------|---------------------------------------|-------|--------|--|
| Rigid<br>(Euclidean) | $[R \ t]_{2 \times 3}^{(3 \times 4)}$ | 3 (6) | Length |  |
|----------------------|---------------------------------------|-------|--------|--|

Similarity

|            |  |       |       |  |
|------------|--|-------|-------|--|
| Similarity | $[SR \ t]_{2 \times 3}^{(3 \times 4)}$ | 4 (7) | Angle |  |
|------------|--|-------|-------|--|

Affine

|        |                                   |        |             |  |
|--------|-----------------------------------|--------|-------------|--|
| Affine | $[A]_{2 \times 3}^{(3 \times 4)}$ | 6 (12) | Parallelism |  |
|--------|-----------------------------------|--------|-------------|--|

Projective

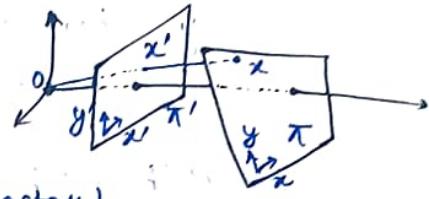
|            |   |        |                |  |
|------------|---|--------|----------------|--|
| Projective | $[\tilde{H}]_{3 \times 3}^{(4 \times 9)}$<br>for 2D<br>for 3D | 8 (15) | Straight Lines |  |
|------------|---|--------|----------------|--|

# ESTIMATING HOMOGRAPHY

## Homography (Plane-to-plane Mapping)

↳ Central projection maps points on one plane to points on another plane, and represented by a linear mapping of homogeneous coordinates  $\mathbf{x}' = \mathbf{H}\mathbf{x}$ .

$$\mathbf{H} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \quad (\text{up to a scalar factor})$$



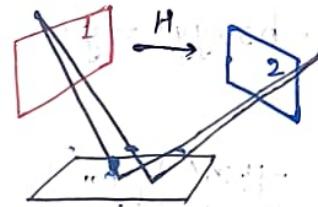
↳ 9 parameters

↳ 8 DoFs

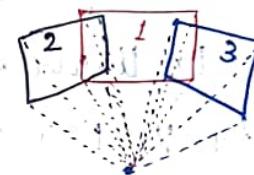
→ There are 2 kinds of homography expressed by the same form:

i) Two-view homography

ii) Homography induced by central projection.



Two-view Homography



Homography induced by central projection

Projectivity: An invertible mapping  $h$  from  $P^2$  to itself such that three points  $x_1, x_2, x_3$  lie on the same line if and only if  $h(x_1), h(x_2), h(x_3)$  do.

↳ = collinearity

= projective transformation

= homography

Setup homography equations:

$$\begin{bmatrix} x' \\ y' \\ w' \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad \text{OR} \quad \begin{bmatrix} x'_1 \\ x'_2 \\ x'_3 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\Rightarrow x' = h_{11}x + h_{12}y + h_{13}$$

$$y' = h_{21}x + h_{22}y + h_{23}$$

$$w' = h_{31}x + h_{32}y + h_{33}$$

- Select 4 points in a plane with known coordinates.
- Form equations.

Converting back to cartesian coordinates,

$$x' = \frac{h_{11}x + h_{12}y + h_{13}}{h_{31}x + h_{32}y + h_{33}}$$

$$y' = \frac{h_{21}x + h_{22}y + h_{23}}{h_{31}x + h_{32}y + h_{33}}$$

(Divide by  $w'$ )

$$\text{OR, } x' = x_1'/x_3' \\ y' = y_1'/x_3'$$

$$\Rightarrow x'(h_{31}x + h_{32}y + h_{33}) = h_{11}x + h_{12}y + h_{13}$$

$$y'(h_{31}x + h_{32}y + h_{33}) = h_{21}x + h_{22}y + h_{23}$$

Set  $h_{33} = 1$  to solve them.

Remark: No calibration at all necessary;

Does not work if  $h_{33} = 0$  in H

$$\Rightarrow h_{11}x + h_{12}y + h_{13} - h_{31}x x' - h_{32}y x' = x'$$

$$h_{21}x + h_{22}y + h_{23} - h_{31}x y' - h_{32}y y' = y'$$

Convert to linear system.

$$\begin{bmatrix} x_1 & y_1 & 1 & 0 & 0 & 0 & -x_1 x_1' & -y_1 x_1' \\ 0 & 0 & 0 & x_1 & y_1 & 1 & -x_1 y_1' & -y_1 y_1' \\ x_2 & y_2 & 1 & 0 & 0 & 0 & -x_2 x_2' & -y_2 x_2' \\ 0 & 0 & 0 & x_2 & y_2 & 1 & -x_2 y_2' & -y_2 y_2' \\ \vdots & \vdots \\ x_n & y_n & 1 & 0 & 0 & 0 & -x_n x_n' & -y_n x_n' \\ 0 & 0 & 0 & x_n & y_n & 1 & -x_n y_n' & -y_n y_n' \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \\ h_{21} \\ h_{22} \\ h_{23} \\ h_{31} \\ h_{32} \end{bmatrix} = \begin{bmatrix} x_1' \\ y_1' \\ x_2' \\ y_2' \\ \vdots \\ x_n' \\ y_n' \end{bmatrix}$$

$\Rightarrow \boxed{Ah = b}$ ,  $A \rightarrow 2N \times 8$  matrix of source & dest. points.  
 $h \rightarrow 8$  element homography vector  
 $b \rightarrow 2N$ -dim. vector of transformed coordinates.

One point pair provides 2 constraints,  
 $8 \text{ DOF} \Rightarrow 4 \text{ point correspondences needed.}$

## Solving H:

Assume  $N (\geq 4)$  point correspondences are given:

$$\{(x_i, y_i), (x'_i, y'_i)\} \quad i=1, \dots, N,$$

$\lambda_i$  is a point-dependent scalar factor.

$$\lambda_i \begin{bmatrix} h_1 & h_2 & h_3 \\ h_4 & h_5 & h_6 \\ h_7 & h_8 & h_9 \end{bmatrix} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} = \begin{bmatrix} x'_i \\ y'_i \\ 1 \end{bmatrix}$$

$$\Rightarrow \lambda_i(h_1 x_i + h_2 y_i + h_3) = x'_i \quad \text{--- (1)}$$

$$\lambda_i(h_4 x_i + h_5 y_i + h_6) = y'_i \quad \text{--- (2)}$$

$$\lambda_i(h_7 x_i + h_8 y_i + h_9) = 1 \quad \text{--- (3)}$$

$$\Rightarrow \lambda_i = 1/(h_7 x_i + h_8 y_i + h_9) \quad [\text{from (3)}] \quad \text{--- (4)}$$

Substitute  $\lambda_i$  in (1), (2),

$$\Rightarrow h_1 x_i + h_2 y_i + h_3 = x'_i (h_7 x_i + h_8 y_i + h_9)$$

$$h_4 x_i + h_5 y_i + h_6 = y'_i (h_7 x_i + h_8 y_i + h_9)$$

In matrix form, we get the following equations:

$$\begin{bmatrix} x_i & y_i & 1 & 0 & 0 & 0 & -x'_i x_i & -y'_i y_i & -x'_i \\ 0 & 0 & 0 & x_i & y_i & 1 & -y'_i x_i & -y'_i y_i & -y'_i \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \\ h_6 \\ h_7 \\ h_8 \\ h_9 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

[ $N$  points  $\equiv 2N$  equations]

$$\Rightarrow Ah = 0 \quad \text{--- (5)}$$

Because  $h$  is homogeneous, the solutions  $h$  and  $\alpha h$  make no difference for any  $\alpha \neq 0$ .

(5) have  $2N$  eqns with 9 unknowns, exact soln cannot always exist.

Thus, solve the least-squared-error solution:

$$\min_h \|Ah\|^2, \text{ such that } \|h\|^2 = 1.$$

$$\text{or, } \arg \min_h \frac{\|Ah\|^2}{\|h\|^2} = \frac{h^T A^T A h}{\|h\|^2} \quad \text{--- (6)}$$

According to the principle of Rayleigh quotient, the soln of (6) is the eigenvector corresponding to the smallest eigenvalue of matrix  $A^T A$ .

Thus, the  $s_{0m}$  can be obtained by performing the eigen-decomposition,  $A^T A = P D P^T$ , where  $D \in \mathbb{R}^{N \times N}$  is a diagonal matrix consisting of eigenvalues (in the g. from large to small) of matrix  $A^T A$ .

The  $s_{0m}$  is the last column of  $P$  (i.e., the eigenvector corresponding to the smallest eigenvalue of matrix  $A^T A$ ).

↳ can be computed by performing singular value decomposition (SVD),  $A = U \Sigma V^T$

(if we arrange the singular values are from large to small (descending order) in the matrix  $\Sigma$ ).

↳ The  $s_{0m}$  is then the last column of  $V$ .

### Singular Value Decomposition

↳ Factorization of a real or complex matrix into a rotation, followed by a scaling followed by another rotation.

↳ It generalizes the eigendecomposition of a square normal matrix with an orthonormal eigenbasis to any  $m \times n$  matrix.

### Applications:

i) Image compression: Reduces interpixel redundancy.

ii) Reduce data from 3D to 2D - dimensionality reduction.

iii) Projecting data to different planes - linear transformations.

### Factorization:

$$A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T$$

$U$ :  $m \times m$  orthogonal matrix

$V$ :  $n \times n$  orthogonal matrix

$\Sigma$ :  $m \times n$  with non-zero element along the diagonal.

- The singular values appear along the main diagonal of  $\Sigma$ .

- They appear in descending order:  $\sigma_1 > \dots > \sigma_n > 0$

where,  $\sigma_1^2 > \dots > \sigma_n^2 > 0$  are eigenvalues of  $A^T A$  &  $A A^T$ .

Assuming  $m \geq n$ :

To factorize  $A$ , consider  $A^T A$  and  $A A^T$ .

Properties:  $ATA$  and  $AAT$  are symmetric matrices.

$$ATA = (ATA)^T = A^T(A^T)^T = ATA$$

SVD:

- calculate  $AAT$ ,  $ATA$ .
- find eigenvalues and  $\Sigma$ .
- calculate  $U$  and  $V$ .
- find the SVD factorization form —  $U\Sigma V^T$ .

Eg  $A = \begin{bmatrix} 2 & 2 \\ -1 & 1 \end{bmatrix} \Rightarrow A^T = \begin{bmatrix} 2 & -1 \\ 2 & 1 \end{bmatrix}$

Step-①: calculate  $AAT$  and  $ATA$ .

$$AAT = \begin{pmatrix} 2 & 2 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} 2 & -1 \\ 2 & 1 \end{pmatrix} = \begin{pmatrix} 8 & 0 \\ 0 & 2 \end{pmatrix}$$

$$ATA = \begin{pmatrix} 2 & -1 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 2 & 2 \\ -1 & 1 \end{pmatrix} = \begin{pmatrix} 5 & 3 \\ 3 & 5 \end{pmatrix}$$

Step-②: Eigenvalues and  $\Sigma$ .

$$|AAT - \lambda I| = 0, |ATA - \lambda I| = 0$$

$$\Leftrightarrow \left| \begin{pmatrix} 8 & 0 \\ 0 & 2 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right| = 0$$

$$\Rightarrow \begin{vmatrix} 8-\lambda & 0 \\ 0 & 2-\lambda \end{vmatrix} = 0$$

$$\Rightarrow (8-\lambda)(2-\lambda) = 0$$

$$\Rightarrow \lambda_1 = 8, \lambda_2 = 2$$

∴ singular values:  $\sigma_1 = \sqrt{\lambda_1} = \sqrt{8}$

$$\sigma_2 = \sqrt{\lambda_2} = \sqrt{2}$$

Therefore,  $\Sigma = \begin{bmatrix} \sqrt{8} & 0 \\ 0 & \sqrt{2} \end{bmatrix}$

Step-③: Finding  $U$ .

$$(AAT - \lambda I)x = 0$$

$$\Rightarrow \begin{pmatrix} 8 & 0 \\ 0 & 2 \end{pmatrix} - \lambda_1 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} x_1 = 0, \begin{pmatrix} 8 & 0 \\ 0 & 2 \end{pmatrix} - \lambda_2 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} x_2 = 0$$

$$\Rightarrow \begin{pmatrix} 8-\lambda_1 & 0 \\ 0 & 2-\lambda_1 \end{pmatrix} x_1 = 0 \quad \Rightarrow \begin{pmatrix} 8-\lambda_2 & 0 \\ 0 & 2-\lambda_2 \end{pmatrix} x_2 = 0$$

$$\Rightarrow \begin{pmatrix} 8-8 & 0 \\ 0 & 2-8 \end{pmatrix} x_1 = 0 \quad \Rightarrow \begin{pmatrix} 8-2 & 0 \\ 0 & 2-2 \end{pmatrix} x_2 = 0$$

$$\Rightarrow \begin{pmatrix} 0 & 0 \\ 0 & -6 \end{pmatrix} x_1 = 0 \quad \Rightarrow \begin{pmatrix} 6 & 0 \\ 0 & 0 \end{pmatrix} x_2 = 0$$

$$\Rightarrow x_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \Rightarrow x_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\Rightarrow \mathbf{x}_1 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\Rightarrow \mathbf{u}_1 = \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|}, \mathbf{u}_2 = \frac{\mathbf{x}_2}{\|\mathbf{x}_2\|}$$

$$\Rightarrow \mathbf{u}_1 = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\therefore \mathbf{U} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$$

Step-④: Finding V.

$$\text{Similarly, } (\mathbf{A}^T \mathbf{A} - \lambda \mathbf{I}) \mathbf{x} = 0$$

$$\Rightarrow \mathbf{v}_1 = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} -1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

$$\therefore \mathbf{V} = \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

Step-⑤: Complete SVD

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$$

$$\begin{bmatrix} 2 & 2 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{8} & 0 \\ 0 & \sqrt{2} \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

SVD compression:

- SVD can compress any form of data.
- SVD takes a matrix (square or non-square) and divides it into 2 orthogonal and a diagonal matrix as a sum of much simpler rank one matrices.
- This allows to approximate the input matrix.
- This new factorized matrices can represent the original matrix with lesser space.

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$$

$$= [\mathbf{u}_1 \dots \mathbf{u}_n] \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} [\mathbf{v}_1^T \dots \mathbf{v}_n^T]^T$$

$$= u_1 \sigma_1 v_1^T + \dots + u_n \sigma_n v_n^T$$

$$= \sum_{i=1}^n \sigma_i u_i v_i^T$$

$\sigma_1 \geq \dots \geq \sigma_n$  (first few terms have largest impact on total sum).

⇒ we can approximate  $\mathbf{A}$  by adding only the first few terms.

## Homography Estimation:

- ① Use the known points to get matrix  $A$ .  $[Ah=0]$
- ② Use SVD in numpy.linalg to get matrix  $V$ .
- ③ The last column of  $V$  will be matrix  $h$ .
- ④ Normalize  $hg$  to 1 and reshape to  $3 \times 3$ .
- ⑤ Get matrix  $H$ .

Applications: Image Mosaicing  $\rightarrow$  stitching multiple overlapping images to get a single, larger image with a wider field of view.

## Existence of Projective Homography

### ① Planar scene:

$x_1$  and  $x_2$  is the 3D point  $x$  expressed in  $C_1$  and  $C_2$ , respectively.

$$x_2 = Rx_1 + t.$$

$N = [n_1, n_2, n_3]^T$  is the unit

normal vector representing the plane  $\pi$  w.r.t.  $C_1$  and  $C_2$ .

$d$  is the  $L^2$  distance from plane to  $C_1$ :

$$N^T x_1 = n_1 x + n_2 y + n_3 z = d;$$

$$\Rightarrow \frac{N^T x_1}{d} = t + x_1 \in \pi.$$

Combining two eqns, we get

$$x_2 = \left( R + \frac{t N^T}{d} \right) x_1,$$

since  $\pi_1 x_1 = x_1$  and  $\pi_2 x_2 = x_2$ , we get

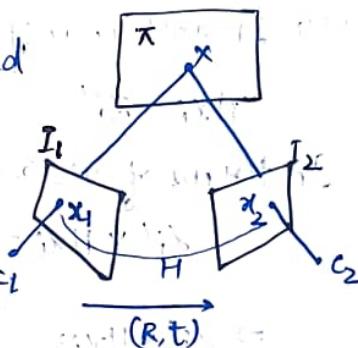
$$x_2 = \underbrace{\left( R + \frac{t N^T}{d} \right)}_{H} x_1$$

② Plane at infinity: Scene is very far away from the camera (e.g. aerial images).

$$H = \left( R + \frac{t N^T}{d} \right) \Rightarrow H_{00} = \lim_{d \rightarrow \infty} \left( R + \frac{t N^T}{d} \right) = R.$$

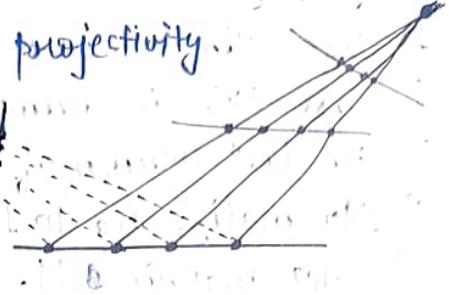
This is a pure rotation, i.e.,  $t = (0, 0, v)^T$ .

$$\Rightarrow H = R.$$



## Projective transformations b/w lines:

- Cross ratio is an invariant under projectivity.
- Cross ratio will have the same value for all the 4 sets of collinear points.



Cross ratio,

$$\text{Cross } (\bar{x}_1, \bar{x}_2; \bar{x}_3, \bar{x}_4) = \frac{|\bar{x}_1 \bar{x}_2| |\bar{x}_3 \bar{x}_4|}{|\bar{x}_1 \bar{x}_3| |\bar{x}_2 \bar{x}_4|}$$

Invariant under  
any projective  
transformation.

$$\text{where } |\bar{x}_i \bar{x}_j| = \det \begin{bmatrix} x_{i1} & x_{j1} \\ x_{i2} & x_{j2} \end{bmatrix}$$

## Direct Linear Transform (DLT):

$$Hx_i = x_i' \quad (\text{Given set of 4 point correspondences: } x_i \leftrightarrow x_i')$$

In terms of vector cross product,

$$x_i' \times Hx_i = 0, \text{ where } Hx_i = \begin{pmatrix} h_1^T x_i \\ h_2^T x_i \\ h_3^T x_i \end{pmatrix}, \quad x_i' = (x_i', y_i', w_i')^T$$

$$\Rightarrow x_i' \times Hx_i = \begin{pmatrix} y_i' h_3^T x_i - w_i h_2^T x_i \\ w_i h_1^T x_i - x_i' h_3^T x_i \\ x_i' h_2^T x_i - y_i h_1^T x_i \end{pmatrix}$$

$$\text{As } h_j^T x_i = x_i^T h_j, \quad j=1, \dots, 3,$$

$$\left\{ \begin{array}{l} \text{only 1st 2 rows are independent} \\ \begin{bmatrix} 0^T & -w_i' x_i^T & y_i' x_i^T \\ w_i x_i^T & 0^T & -x_i' x_i^T \\ -y_i' x_i^T & x_i' x_i^T & 0^T \end{bmatrix} \begin{pmatrix} h_1' \\ h_2 \\ h_3 \end{pmatrix} = 0 \end{array} \right. + (a_1 b_2 - a_2 b_1) R$$

$$\Rightarrow \begin{pmatrix} 0^T & -w_i' x_i^T & y_i' x_i^T \\ w_i x_i^T & 0^T & -x_i' x_i^T \\ 0^T & 0^T & 0^T \end{pmatrix} \begin{pmatrix} h_1' \\ h_2 \\ h_3 \end{pmatrix} = 0$$

$$\Rightarrow A_i h = 0$$

[The 3rd row is obtained, upto scale, from sum of  $x_i'$  times 1st row and  $y_i$  times the second.]

$A_i \rightarrow 2 \times 9$  matrix

$h \rightarrow$  a 9-vector made up of all elements in  $H$ , i.e.,

$$h = \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix}, \quad H = \begin{pmatrix} h_1 & h_2 & h_3 \\ h_4 & h_5 & h_6 \\ h_7 & h_8 & h_9 \end{pmatrix}, \quad \text{with } h_i \text{ the } i^{\text{th}}$$

element of  $h$ ;  $w_i$  chosen]

## Least square solutions:

- ↳ Point correspondences are corrupted with noise.
- ↳ An exact soln for  $Ah=0$  does not exist!
- minimize  $\|Ah\|$  subject to  $\|h\|=1$ .

↳ Obtained by taking the g-vector right null-space

Right null-space: Right singular value that corresponds to the smallest singular value ( $\sigma_g$ ) in the SVD of  $A$  ( $v_g$ ),

$$\text{SVD}(A) = [U_1, U_2, \dots, U_r] \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{bmatrix} [V_1, V_2, \dots, V_r]^T$$

left singular  
vector ( $r \times r$ )      singular values      right singular  
vector ( $g \times g$ )

$$m > n, \text{rank}(A) = r.$$

$\sigma_{n-r}, \dots, \sigma_n = 0$ , i.e.,  $\text{rank}(A) = r$ , if  $A$  is not corrupted by noise and an exact soln for  $Ah=0$  exists!

If  $A$  is noise corrupted,  $\sigma_{n-r}, \dots, \sigma_n \neq 0$ .

$$A = U\Sigma V^T \Rightarrow Av = U\Sigma V^T v = U\Sigma v_i \sigma_i$$

$\|Av\|$  is minimized when  $v_i \sigma_i$  is at its min<sup>n</sup>, i.e., smallest singular value, i.e.,  $\sigma_n$ .

The soln of  $\arg\min_h \|Ah\|$  s.t.  $\|h\|=1$

is given by setting  $h = v_n$

(last column of  $V$ ).

Degeneracy in homography estimation:  $[Ah=0]$

↳ Rank of  $A$  drops below 8 if 3 of the min<sup>n</sup> 4 points correspondences are collinear.

↳ cannot solve for  $h \Rightarrow$  critical configuration or degeneracy.

→ A projective transformation must preserve collinearity.

Data Normalization: Carried out by a transformation of points:

- Points are translated so that their centroid is at origin.
- Points are then scaled so that the average distance from origin is equal to  $\sqrt{2}$ .
- Transformation is then applied to each of 2 images independently
  - ↳ Average point is equal to  $(1, 1, 1)^T$  after normalization.
  - ↳ No magnitude difference in linear eqn  $Ah=0$ .

→ Data normalization is an essential step in DLT algorithm.

↳ Must not be considered optional.

### Distances in DLT:

Algebraic distance:

DLT minimizes  $\|Ah\|$ .

$e = Ah \rightarrow$  residual error.

$e_i \rightarrow$  partial error for each  $x_i \rightarrow x'_i$

algebraic distance:

$$\text{alg.}(x'_i, Hx_i)^2 = \|e_i\|^2 = \left\| \begin{bmatrix} 0^T & -w_i^T x_i^T & -y_i^T x_i^T \\ -w_i^T x_i^T & 0^T & -x_i^T x_i^T \end{bmatrix} h \right\|^2$$

$$\text{alg.}(x_1, x_2)^2 = a_1^2 + a_2^2, \text{ where } a = (a_1, a_2, a_3)^T = x_1 x_2$$

$$\sum \text{alg.}(x_i, Hx_i)^2 = \sum \|e_i\|^2 = \|Ah\|^2 = \|e\|^2.$$

Geometric distance:

Error in one image,

$$\hat{H} = \underset{H}{\operatorname{argmin}} \sum_i d(x_i, H\bar{x}_i)^2$$

Symmetric transfer error, (minimum  $|s_i - \hat{s}_i|$  (in image))

$$\hat{H} = \underset{H}{\operatorname{argmin}} \sum_i d(x_i, H^T x_i)^2 + d(x'_i, H\bar{x}_i)^2$$

Reprojection error,

$$(\hat{H}, \hat{x}_i, \hat{x}'_i) = \underset{H, \hat{x}, \hat{x}'}{\operatorname{argmin}} \sum_i d(x_i, \hat{x}_i)^2 + d(x'_i, \hat{x}'_i)^2, \text{ subject to } \hat{x}'_i = \hat{H}\hat{x}_i.$$

Transfer error,

$$\sum_i d(x'_i, Hx_i)^2$$

Campson error:

$C_H(\bar{x}) = 0$ : cost fn  $Ah = 0$  satisfied by pt.  $x = (x, y, x', y')^T$

$\hat{x}$ : desired pt. such that  $C_H(\hat{x}) = 0$ , where  $\delta x = \hat{x} - x$ .

By Taylor expansion,

$$C_H(x + \delta x) = C_H(x) + \left( \frac{\partial C_H}{\partial x} \right) \delta x = 0.$$

Approximated cost fn:  $J \delta x = -\epsilon$ ,

$J$ : partial-derivative matrix

$\epsilon$ : cost  $C_H(x)$  associated with  $x$ .

## Camera Calibration and Pose Estimation (Camera Localization)

## Camera Calibration

↳ Process to determine

↳ extrinsic parameters ( $R, T$ ) of a camera.

↳ intrinsic parameters ( $K$  plus lens distortion).  
point in camera frame

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K [R | T] \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix}}_{K: \text{intrinsic/ calibration matrix}} \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix}$$

Extrinsic matrix

pt. in the camera frame

→ Estimating of extrinsic parameters with known intrinsic parameters: camera localization.

→ Neglect lens distortion (eg, barrel distortion) for now.

Mapping a 3-D world point to image point:

$$x = P \cdot X$$

pixel coordinate  $\downarrow$   $\rightarrow$  world coordinate  
formation

Wanted: Extrinsic & intrinsic parameters of a camera.

Given: Coordinates of object points (control points).

Observed: Coordinates of those known 3D object points in the image.

## Coordinate Systems:

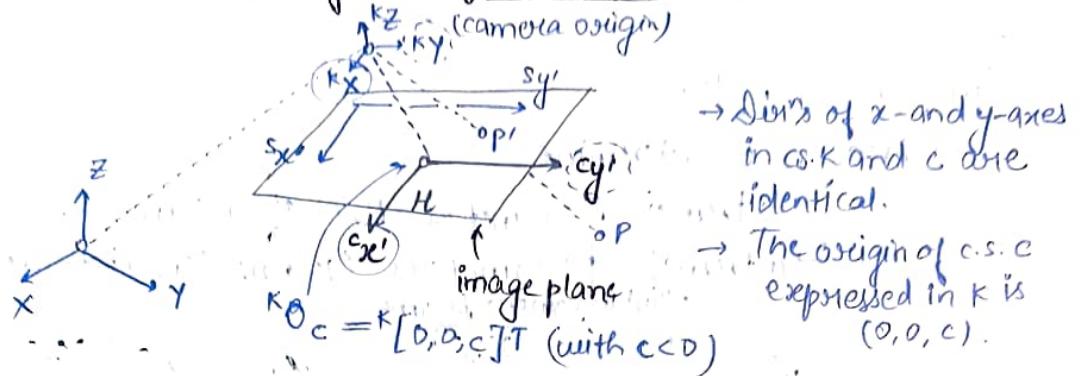
- ① World/Object coordinate system  $S_o$ :  $[x, y, z]^T$  (no index)
  - ② Camera coordinate system  $S_k$ :  $[k_x, k_y, k_z]^T$
  - ③ Image plane coordinate system  $S_c$ :  $[c_x, c_y]^T$
  - ④ Sensor coordinate system  $S_s$ :  $[s_x, s_y]^T$ .

$$\begin{bmatrix} sx \\ sy \\ sz \end{bmatrix} = SH_C \cdot CP_k \cdot KH_o \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

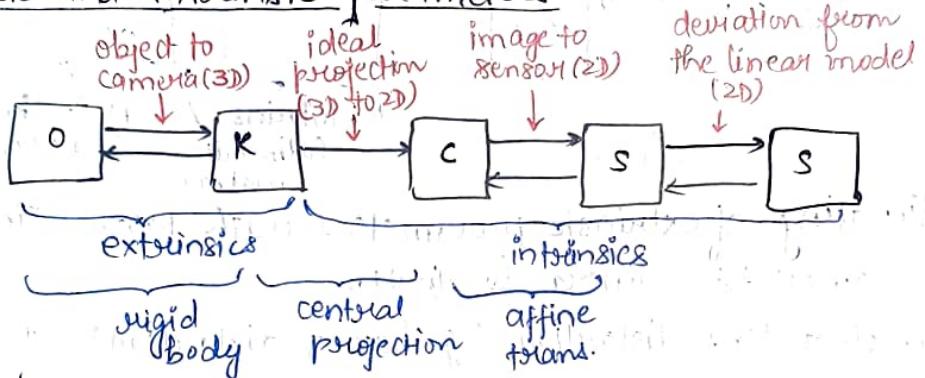
↓      ↓      ↓      ↓

in sensor system      image plane      camera to image      object to camera      in the object system

x and y axis coordinate systems of K and c:



Extrinsic and intrinsic parameters:



Extrinsic parameters: Describe the pose (=position and heading) of the camera w.r.t. world.  
 ↳ Invertible transformations.

↳ 6 parameters: 3 for the position + 3 for heading.

- Point P in world coordinates:  $\mathbf{x}_p = [x_p, y_p, z_p]^T$
- Center O of the projection:  $\mathbf{x}_o = [x_o, y_o, z_o]^T$  (camera coord. origin) (or  $z$  or  $z_o$ )

→ In Euclidean coordinate systems:

- Translation b/w origin of world c.s. and camera c.s.:  
 $\mathbf{x}_o = [x_o, y_o, z_o]^T$

- Rotation R from  $s_o$  to  $s_K$

$$K\mathbf{x}_p = R(\mathbf{x}_p - \mathbf{x}_o)$$

→ In homogeneous c.s.,

$$\begin{bmatrix} K\mathbf{x}_p \\ 1 \end{bmatrix} = \begin{bmatrix} R & 0 \\ 0^T & 1 \end{bmatrix} \begin{bmatrix} I_3 - \mathbf{x}_o \\ 0^T \end{bmatrix} \begin{bmatrix} \mathbf{x}_p \\ 1 \end{bmatrix}$$

Euclidean H.C.

$$= \begin{bmatrix} R & -R\mathbf{x}_o \\ 0^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{x}_p \\ 1 \end{bmatrix}$$

or,  $K\mathbf{x}_p = K\mathbf{H}\mathbf{x}_p$  with  $K\mathbf{H} = \begin{bmatrix} R & -R\mathbf{x}_o \\ 0^T & 1 \end{bmatrix}$

Intrinsic parameters: Projecting points from camera c.s. to sensor c.s.

↳ Invertible: image plane to sensor, model deviations

↳ Non-invertible: central projection.

→ 3 steps mapping:

① Ideal projection (perspective) to the image plane.

② Mapping to the ~~s-~~ system sensor c.s. (where the pixels are).

③ Compensation for the fact that the 2 previous mappings are idealized.

→ Assumptions:

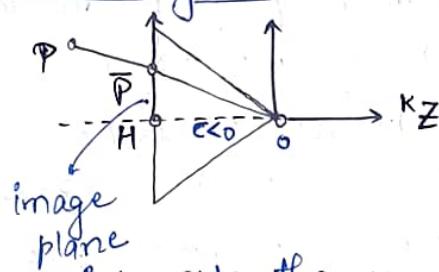
① Distortion-free lens.

② All rays are straight lines and pass through the projection center (origin of the camera c.s.  $\mathbf{s}_k$ ).

③ Focal point and principal point lie on the optical axis.

④ Distance from the camera origin to the image plane is the constant 'c'.

Image coordinate systems:



↳ Most popular image c.s.

(Rotated 180° from physically motivated c.s.)

Through Intercept theorem,

$$c_{x\bar{P}} := k_{x\bar{P}} = c \frac{k_{X_P}}{k_{Z_P}}$$

$$c_{y\bar{P}} := k_{y\bar{P}} = c \frac{k_{X_P}}{k_{Z_P}}$$

$$\left( c = k_{z\bar{P}} = c \frac{k_{Z_P}}{k_{Z_P}} \right)$$

In homogeneous c.s.,

$$\begin{bmatrix} k_{U\bar{P}} \\ k_{V\bar{P}} \\ k_{W\bar{P}} \\ k_{T\bar{P}} \end{bmatrix} = \begin{bmatrix} c & 0 & 0 & 0 \\ 0 & c & 0 & 0 \\ 0 & 0 & c & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} k_{X_P} \\ k_{Y_P} \\ k_{Z_P} \\ 1 \end{bmatrix}$$

↓ Drop the 3rd row

$$\begin{bmatrix} c_{x\bar{P}} \\ c_{y\bar{P}} \\ 1 \end{bmatrix} = \begin{bmatrix} c_{U\bar{P}} \\ c_{V\bar{P}} \\ c_{W\bar{P}} \end{bmatrix} = \begin{bmatrix} c & 0 & 0 & 0 \\ 0 & c & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} k_{X_P} \\ k_{Y_P} \\ k_{Z_P} \\ 1 \end{bmatrix} = \begin{bmatrix} c k_{X_P} \\ c k_{Y_P} \\ c k_{Z_P} \\ 1 \end{bmatrix} = \begin{bmatrix} c \frac{k_{X_P}}{k_{Z_P}} \\ c \frac{k_{Y_P}}{k_{Z_P}} \\ c \frac{k_{Z_P}}{k_{Z_P}} \\ 1 \end{bmatrix}$$

Thus, for any point,

$${}^c x_p = {}^c P_K {}^k x_p \text{, with } {}^c P_K = \begin{bmatrix} C & 0 & 0 \\ 0 & C & 0 \\ 0 & 0 & I \end{bmatrix}$$

This gives,

$${}^c x = {}^c P X \text{, with } {}^c P = {}^c P_K {}^K H = \begin{bmatrix} C & 0 & 0 \\ 0 & C & 0 \\ 0 & 0 & I \end{bmatrix} \begin{bmatrix} R & -R x_0 \\ 0 & I \end{bmatrix}$$

Calibration

matrix

(for ideal camera)

$${}^c K = \begin{bmatrix} C & 0 & 0 \\ 0 & C & 0 \\ 0 & 0 & I \end{bmatrix}$$

∴ Overall mapping:

$${}^c p = {}^c K [R] - R x_0 = {}^c K R [I_3] - x_0$$

coordinates of  ${}^c x$ :

$$\begin{bmatrix} {}^c u' \\ {}^c v' \\ {}^c w' \end{bmatrix} = \begin{bmatrix} C & 0 & 0 \\ 0 & C & 0 \\ 0 & 0 & I \end{bmatrix} \begin{bmatrix} \gamma_{11} & \gamma_{12} & \gamma_{13} \\ \gamma_{21} & \gamma_{22} & \gamma_{23} \\ \gamma_{31} & \gamma_{32} & \gamma_{33} \end{bmatrix} \begin{bmatrix} x - x_0 \\ y - y_0 \\ z - z_0 \end{bmatrix}$$

In Euclidean coordinates,

this leads to collinearity eq<sup>n</sup> for the image coordinates.

$${}^c x = c \frac{\gamma_{11}(x - x_0) + \gamma_{12}(y - y_0) + \gamma_{13}(z - z_0)}{\gamma_{31}(x - x_0) + \gamma_{32}(y - y_0) + \gamma_{33}(z - z_0)}$$

$${}^c y = c \frac{\gamma_{21}(x - x_0) + \gamma_{22}(y - y_0) + \gamma_{23}(z - z_0)}{\gamma_{31}(x - x_0) + \gamma_{32}(y - y_0) + \gamma_{33}(z - z_0)}$$

Mapping to the sensor (without linear errors): Image → sensor

Linear errors:

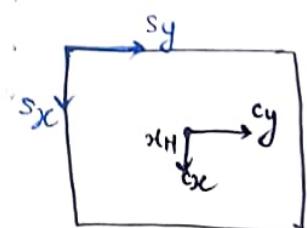
- ① Location of the principal point in the image.
- ② scale difference in  $x$  and  $y$ -based on the chip design.
- ③ shear compensation.

Location to principal point:

↪ The origin of the sensor system is not at the principal pt.

↪ Compensation through a shift.

$${}^S H_c = \begin{bmatrix} 1 & 0 & x_H \\ 0 & 1 & y_H \\ 0 & 0 & 1 \end{bmatrix}$$



## Shear and scale differences:

↳ Scale difference  $m$  in  $x$  and  $y$ .

↳ shear compensation  $s$  (for digital cameras, typically  $s \approx 0$ )

$$\therefore sH_c = \begin{bmatrix} 1 & s & x_H \\ 0 & 1+m & y_H \\ 0 & 0 & 1 \end{bmatrix}$$

$$\text{Hence, } s_x = sH_c \cdot {}^cKR [I_3 | -x_0] x$$

→ Often,  $sH_c$  is combined with calibration matrix  ${}^cK$ .

$$K = sH_c {}^cK$$

$$= \begin{bmatrix} 1 & s & x_H \\ 0 & 1+m & y_H \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} c & 0 & 0 \\ 0 & c & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} c & cs & x_H \\ 0 & c(1+m) & y_H \\ 0 & 0 & 1 \end{bmatrix}$$

↓  
affine transformation  
(K)

↳ 5 parameters:

- camera constant ( $c$ )
- principal point ( $x_H, y_H$ )
- scale difference ( $m$ )
- shear ( $s$ )

→ Affine camera: camera with an affine mapping to the sensor c.s. (after central projection is applied).

## Direct Linear Transform (DLT):

$x = P X$ , with  $P = {}^cKR [I_3 | -x_0]$  → model of  
maps any object pt.  $x$  to the image pt.  $x$ . ↳ 11 parameters: affine camera

- 6 extrinsic ( $R, x_0$ )
- 5 intrinsic ( $c, x_H, y_H, m, s$ )

↳ Homogeneous projection matrix

↳ both in- and extrinsic parameters.

In Euclidean,

$$s_x = \frac{p_{11}X + p_{12}Y + p_{13}Z + p_{14}}{p_{31}X + p_{32}Y + p_{33}Z + p_{34}}$$

$$s_y = \frac{p_{21}X + p_{22}Y + p_{23}Z + p_{24}}{p_{31}X + p_{32}Y + p_{33}Z + p_{34}}$$

DLT maps obj. pt.  $X$  to image pt.  $x$ ,

$$\begin{aligned} x &= \underbrace{\begin{pmatrix} K & R \\ (3 \times 3) & (3 \times 3) \end{pmatrix}}_{(3 \times 4)} \underbrace{\begin{pmatrix} I_3 & -x_0 \\ (3 \times 3) & (3 \times 1) \end{pmatrix}}_{3 \times 4} X \quad \left. \begin{array}{l} K \rightarrow \text{intrinsic} \\ x_0, R \rightarrow \text{extrinsic} \end{array} \right. \\ &= \underbrace{\begin{pmatrix} P \\ (3 \times 4) \end{pmatrix}}_{(4 \times 1)} X \end{aligned}$$

## Computing the Orientation of an Uncalibrated camera:

### Calibrated camera

- ↳ 6 unknowns
- ↳ Need atleast 3 points
- ↳ solved by spatial resection

### Uncalibrated camera

- ↳ 11 unknowns
- ↳ Need atleast 6 points
- ↳ Assuming the model of an affine camera.
- ↳ Problem - solved by DLT.

DLT eqn:

$$x_i = P \quad X_i = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} X_i$$

∴ Define 3 vectors:  $A = \begin{bmatrix} p_{11} \\ p_{12} \\ p_{13} \\ p_{14} \end{bmatrix}$ ,  $B = \begin{bmatrix} p_{21} \\ p_{22} \\ p_{23} \\ p_{24} \end{bmatrix}$ ,  $C = \begin{bmatrix} p_{31} \\ p_{32} \\ p_{33} \\ p_{34} \end{bmatrix}$ .

$$\text{Then, } x_i = P X_i = \begin{bmatrix} A^T \\ B^T \\ C^T \end{bmatrix} X_i = \begin{bmatrix} A^T X_i \\ B^T X_i \\ C^T X_i \end{bmatrix}$$

$$\Rightarrow x_i = \frac{u_i}{w_i} = \frac{A^T X_i}{C^T X_i}, \quad y_i = \frac{v_i}{w_i} = \frac{B^T X_i}{C^T X_i}$$

$$\Rightarrow x_i C^T X_i - A^T X_i = 0$$

$$y_i C^T X_i - B^T X_i = 0$$

$$\Rightarrow \begin{aligned} -x_i^T A + x_i^T C &= 0 \\ -x_i^T B + y_i^T C &= 0 \end{aligned} \quad \left. \begin{array}{l} \text{system of linear eq's} \\ \text{in parameters } A, B, C. \end{array} \right.$$

$$\Rightarrow \begin{aligned} a_{x_i}^T p &= 0 && \text{collect elements of } P \text{ within a} \\ a_{y_i}^T p &= 0 && \text{parameter vector } p \end{aligned}$$

↙ rows of  $P$  as column-vectors

$$a_{x_i}^T = (x_i^T, 0^T, x_i x_i^T)$$

$$a_{y_i}^T = (0^T, -x_i^T, y_i x_i^T)$$

homogeneous system

↳ solve using SVD (finding null space of  $A$ ).

↳ 'No sol'n' if all  $x_i$  are on

• if all  $x_i$  and projection center  $x_0$  are located on a twisted cubic curve.

## Decomposition of P:

$$P = [KR| -KRX_0] = [H|h]$$

with  $H = KR$ ,  $h = -KRX_0$

$K \rightarrow$  triangular matrix  
 $R \rightarrow$  rotation matrix

we get ~~decom~~ projection center through

$$X_0 = -H^{-1}h.$$

$H = KR$  structure: QR-decomposition

Eg.  $A = \begin{pmatrix} 1 & q_1 & q_2 & q_3 \\ 0 & \|q_1\| & q_2^\perp & q_3^\perp \\ 0 & 0 & \|q_2^\perp\| & q_3^\perp \\ 0 & 0 & 0 & \|q_3^\perp\| \end{pmatrix}$   $\Rightarrow Q = \begin{pmatrix} 1 & q_1 & q_2 & q_3 \\ 0 & 1 & q_2^\perp & q_3^\perp \\ 0 & 0 & 1 & q_3^\perp \\ 0 & 0 & 0 & 1 \end{pmatrix}$

- $q_1 = \frac{a_1}{\|a_1\|}$  ( $\because \|q_1\| = 1$ )
- $a_2^\perp = a_2 - \langle a_2, q_1 \rangle q_1$
- $q_2 = \frac{a_2^\perp}{\|a_2^\perp\|}$
- $a_3^\perp = a_3 - \langle a_3, q_1 \rangle q_1 - \langle a_3, q_2 \rangle q_2$
- $q_3 = \frac{a_3^\perp}{\|a_3^\perp\|}$

$$\therefore H^{-1} = (KR)^{-1} = R^{-1}K^{-1} = R^T K^{-1}$$

$\rightarrow H = KR \rightarrow$  homogeneous matrix  $\Rightarrow$  calibration matrix

Normalize  $K \leftarrow \frac{1}{K_{33}} K$

$\rightarrow$  To get -ve camera constant,

choose  $K \leftarrow KR(3, \pi)$ ,  $R \leftarrow R(3, \pi)R$

using  $R(3, \pi) = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{pmatrix}$

## Camera Calibration by Zhang's Method

$$x = KR [I_3 | -X_0] X$$

compute 5 intrinsic parameters.

observed image point  $c, s, m,$   
 $x_1, y_1$   $\uparrow$  3 translations  
 $\uparrow$  3 rotations  $\uparrow$  control point (given)

Using 2D checkerboard

$\hookrightarrow$  Observed 2D pattern

$\hookrightarrow$  Known size and structure

$\hookrightarrow$  Set the world coordinate system to the corner of the checkerboard.

$\hookrightarrow$  All points on the checkerboard lie in the XY plane, i.e.,  $Z=0$ .

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} c & c_8 & x_H \\ 0 & c(1+m) & y_H \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mathcal{H}_{11} & \mathcal{H}_{12} & \cancel{\mathcal{H}_{13}} & t_1 \\ \mathcal{H}_{21} & \mathcal{H}_{22} & \cancel{\mathcal{H}_{23}} & t_2 \\ \mathcal{H}_{31} & \mathcal{H}_{32} & \cancel{\mathcal{H}_{33}} & t_3 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

→ Z coordinate of each pt. on the checkerboard  
Delete 3rd column of extrinsic matrx.

$\Downarrow$

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} c & c_8 & x_H \\ 0 & c(1+m) & y_H \\ 0 & 0 & 1 \end{bmatrix}}_{K} \underbrace{\begin{bmatrix} \mathcal{H}_{11} & \mathcal{H}_{12} & t_1 \\ \mathcal{H}_{21} & \mathcal{H}_{22} & t_2 \\ \mathcal{H}_{31} & \mathcal{H}_{32} & t_3 \end{bmatrix}}_{[\mathcal{H}_1, \mathcal{H}_2, t]} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$H = [h_1, h_2, h_3]$

For multiple observed pts.:  $\begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} = H_{3 \times 3} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix}, i=1, \dots, I$

- Estimate  $3 \times 3$  homography instead of  $3 \times 4$  projection matrix.
- Solve to estimate  $H$ . → Need at least 4 points
- No direct decomposition as in DLT. ( $H \rightarrow 8 \text{ DOF}$ ) → 2 observations from each point.
- Exploit constraints of the parameter.

$$[\mathcal{H}_1, \mathcal{H}_2, t] = K^{-1} [h_1, h_2, h_3]$$

$$\Rightarrow \mathcal{H}_1 = K^{-1} h_1, \mathcal{H}_2 = K^{-1} h_2.$$

$\mathcal{H}_1, \mathcal{H}_2, \mathcal{H}_3$  form an orthonormal basis

$$\Rightarrow \mathcal{H}_1^T \mathcal{H}_2 = 0, \quad \|\mathcal{H}_1\| = \|\mathcal{H}_2\| = 1.$$

$$\mathcal{H}_1^T \mathcal{H}_2 = 0 \Rightarrow h_1^T K^{-T} K^{-1} h_2 = 0$$

$$\|\mathcal{H}_1\| = \|\mathcal{H}_2\| = 1 \Rightarrow h_1^T K^{-T} K^{-1} h_1 = h_2^T K^{-T} K^{-1} h_2$$

$$\Rightarrow h_1^T K^{-T} K^{-1} h_1 - h_2^T K^{-T} K^{-1} h_2 = 0$$

Define symmetric and toe definite matrix:  $B = K^{-T} K^{-1}$

$$h_1^T B h_2 = 0$$

$$h_1^T B h_1 - h_2^T B h_2 = 0$$

From B, the calibration matrix can be recovered through Cholesky decomposition.

$$B = \begin{pmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{pmatrix}$$

$$\text{chol}(B) = A A^T$$

$$\Rightarrow A = K^{-T}$$

→ If we know B, we can compute K.

↪ B: unknown

↪ h: known

↪ 2 equations relating B and h.

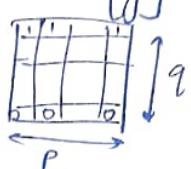
Cholesky decomposition:

$$l_{ki} = a_{ki} - \sum_{j=1}^{i-1} l_{ij} l_{kj}, \quad l_{kk} = \sqrt{a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2}$$

# Image Processing and Feature Extraction

## Image

- ↳ (gray scale)  $\rightarrow$  2D function ; domain  $x = \begin{bmatrix} x \\ y \end{bmatrix}$
- ↳ Types : ① Binary  $\rightarrow$ 
  - 0: Black
  - 1: White
- ② Gray scale
- ③ Color



## Image Transformations

- ↳ Filtering : changes pixel values :  $G(x) = h \{ F(x) \}$  original image
- ↳ Warping : changes pixel locations :  $G(x) = F(h\{x\})$  changes range
- ↳ changes domain of image function

### Image filtering

- ↳ Neighborhood operation (filtering)

### Point operation (point processing)

- original:  $x$
- Darken:  $x - 128$
- Lower contrast:  $x/2$
- Non-linear lower contrast:  $(\frac{x}{255})^k \times 255$
- Invert:  $255 - x$
- Lighten:  $x + 128$
- Raise contrast:  $x \times 2$
- Non-linear raise contrast:  $(\frac{x}{255})^2 \times 255$

→ Images are sampled.

### Resolution: Sampling parameter

- ↳ defined in dots per inch (DPI) ; standard value: 72 dpi
- ↳ equivalent measure of spatial pixel density.

### Images are sampled and quantized:

- ↳ An image contains discrete no. of pixels.

Eg. Pixel value : gray scale (or intensity) :  $[0, 255]$

color : RGB :  $[R, G, B]$

Lab :  $[L, a, b]$

HSV :  $[H, S, V]$

- ↳ Errors due to sampling and quantization.

## Histogram: A type of image.

- ↳ Provides frequency of brightness (intensity) value in the image.
- ↳ Captures the distribution of gray levels in the image (how frequent).



## Image as Discrete Functions: Digital (discrete)

- ↳ Sample the 2D space on a regular grid.

- ↳ Represented as a matrix of integer values in cartesian coordinates:

$$f[n, m] = \begin{bmatrix} f[-1, -1] & f[0, -1] & f[1, -1] \\ f[-1, 0] & f[0, 0] & f[1, 0] \\ f[-1, 1] & f[0, 1] & f[1, 1] \end{bmatrix}$$

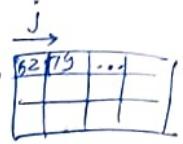


Image:  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^M$

$f(x, y)$  → gives intensity at position  $(x, y)$ .

- ↳ Defined over a rectangle with a finite range:

$$f: \underbrace{[a, b] \times [c, d]}_{\text{domain support}} \rightarrow \underbrace{[0, 255]}_{\text{range}}$$

$$\text{color image: } f(x, y) = \begin{bmatrix} g_1(x, y) \\ g_2(x, y) \\ g_3(x, y) \end{bmatrix}$$

## Linear shift-invariant image filtering:

- ↳ Replace each pixel by a linear combination of its neighbors (and possibly itself).

- ↳ The combination is determined by filter's kernel.

## Systems and Filters

shifted to all pixel locations

Filtering → Form a new image whose pixel values are transformed from original pixel values.

- ↳ To extract useful info from images, or transform images into another domain where we can modify/ enhance image properties:

↳ Features (edges, corners, blobs)

↳ Super-resolution, in-painting, de-noising.

$f[n, m] \xrightarrow{\text{(pimage)}} \boxed{\text{Systems}} \xrightarrow{\text{(filter)}} g[n, m] \xrightarrow{\text{(o/p image)}} (n, m): \begin{cases} \text{spatial position in the image} \end{cases}$

Properties:  $\text{filter}(\text{im}, f_1 + f_2) = \text{filter}(\text{im}, f_1) + \text{filter}(\text{im}, f_2)$

$$C^* \text{filter}(\text{im}, f_1) = \text{filter}(\text{im}, C^* f_1).$$

Box Filter: (or 2D mean filter or square mean filter)

Kernel:  $g[\cdot, \cdot] = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$

↳ Replaces pixel with local average

↳ Has smoothening(blurring) effect. (approximate Gaussian blur.)

Convolution:

1D continuous signals:  $(f * g)(x) = \int_{-\infty}^{\infty} f(y) g(x-y) dy$

2D discrete signals:  $(f * g)(x, y) = \sum_{i,j=-\infty}^{(\infty)} f(i, j) I(x-i, y-j)$

if  $f(i, j) \neq 0$  within only  $-1 \leq i, j \leq 1$ .

Correlation:

$$(f * g)(x, y) = \sum_{i,j=-\infty}^{\infty} f(i, j) I(x+i, y+i)$$

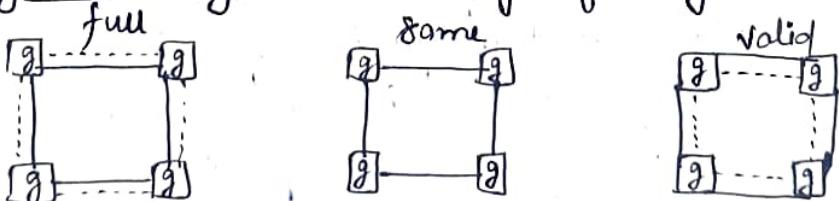
[Most kernels are symmetric.]

separable filters: If it can be written as a product of a column and a row.

E.g. box filter:  $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$

→ 2D convolution with a separable filter is equivalent to two 1D convolutions.

Handling boundary issues in image filtering:



Border padding

- ↳ Zero pad
- ↳ Circular
- ↳ Replicate
- ↳ Symmetric

Moving Average filter:  $\rightarrow$  shift-invariant!

$$g[n, m] = \frac{1}{9} \sum_{k=-1}^{m+1} \sum_{l=-1}^{m+1} f[k, l]$$

$$= \frac{1}{9} \sum_{k=-1}^1 \sum_{l=-1}^1 f[n-k, m-l]$$

$$\frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

$\hookrightarrow$  Smoothing effect (remove sharp features).

Image segmentation: Based on a simple threshold.  $\rightarrow$  Not linear

$$g[n, m] = \begin{cases} 255, & f[n, m] > 100 \\ 0, & \text{otherwise.} \end{cases}$$

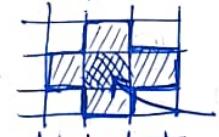
### Properties of Systems

- Linear system
- ① Additivity:  $S[f_i[n, m] + f_j[n, m]] = S[f_i[n, m]] + S[f_j[n, m]]$
  - ② Homogeneity:  $S[\alpha f_i[n, m]] = \alpha S[f_i[n, m]]$
  - ③ Superposition:  $S[\alpha f_i[n, m] + \beta f_j[n, m]] = \alpha S[f_i[n, m]] + \beta S[f_j[n, m]]$
  - ④ Stability:  $|f[n, m]| \leq K \Rightarrow |g[n, m]| \leq cK$
  - ⑤ Invertibility:  $S^{-1}[S[f[n, m]]] = f[n, m]$ .
  - ⑥ Causality: For  $n < n_0, m < m_0$ , if  $f[n, m] = 0 \Rightarrow g[n, m] = 0$
  - shift invariant system
    - ⑦ Shift Invariance:  $f[n-n_0, m-m_0] \xrightarrow{S} g[n-n_0, m-m_0]$

### Causality in Images:

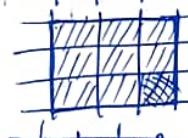
supposed

Non-causal:

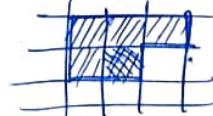


o/p pixel

causal:



Semi-causal:



### Characterize a Linear Shift Invariant (LSI) System

2D impulse function:

$\hookrightarrow 1$  at  $[0, 1]$

$\hookrightarrow 0$  everywhere else

$$\delta_2[m, n]$$

Impulse Response to moving filter:

$$\delta_2 \xrightarrow{S} h[n, m] = \frac{1}{9} \sum_{k=-1}^1 \sum_{l=-1}^1 \delta_2[n-k, m-l]$$

$$= \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

General LSI system:  $f[n, m] = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} f[k, l] \times s_2[n-k, m-l]$

## 2D Convolution:

$$f[n, m] * h[n, m] = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} f[k, l] h[n-k, m-l] \rightarrow \text{Filtering operation.}$$

Eg. Kernel:

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

→ Filtered image  
↓  
no change

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

→ shifted weight  
by 1 pixel

$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

→ Blur

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$-\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

→ Sharpening filter

→ Accentuates  
differences with  
local average

Image Support and Edge Effect: Finite support signals convolution

$$\begin{bmatrix} \quad & \quad & \quad \\ \quad & \quad & \quad \\ N_1 \times M_1 & \end{bmatrix}$$

$$* \begin{bmatrix} \quad & \quad & \quad \\ \quad & \quad & \quad \\ N_2 \times M_2 & \end{bmatrix}$$

$$\begin{bmatrix} \quad & \quad & \quad \\ \quad & \quad & \quad \\ N_1 + N_2 - 1 \times M_1 + M_2 - 1 & \end{bmatrix}$$

## Convolution as Matrix Multiplication:

① Define I/p and filter

$$I = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} (m_1 \times n_1)$$

$$F = \begin{bmatrix} 10 & 20 \\ 30 & 40 \end{bmatrix} (m_2 \times n_2)$$

② calculate final o/p size:

$$(m_1 + m_2 - 1) \times (n_1 + n_2 - 1)$$

③ zero-pad filter matrix

$$\text{Zero-padded } F = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 10 & 20 & 0 & 0 \\ 30 & 40 & 0 & 0 \end{bmatrix}$$

④ Create Toeplitz matrix for each row of zero-padded filter.

$$F_0 = \begin{bmatrix} 30 & 0 & 0 \\ 40 & 30 & 0 \\ 0 & 40 & 30 \\ 0 & 0 & 40 \end{bmatrix}$$

$$F_1 = \begin{bmatrix} 10 & 0 & 0 \\ 20 & 10 & 0 \\ 0 & 20 & 10 \\ 0 & 0 & 20 \end{bmatrix}$$

$$F_2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

⑤ Create a doubly blocked Toeplitz matrix

$$\text{doubly blocked} = \begin{bmatrix} F_0 & 0 \\ F_1 & F_0 \\ F_2 & F_1 \end{bmatrix}$$

$$= \begin{bmatrix} 30 & 0 & 0 & 0 & 0 \\ 40 & 30 & 0 & 0 & 0 \\ 0 & 40 & 30 & 0 & 0 \\ 0 & 0 & 40 & 30 & 0 \\ 10 & 0 & 0 & 40 & 30 \\ 20 & 10 & 0 & 0 & 40 & 30 \\ 0 & 20 & 10 & 0 & 0 & 40 & 30 \\ 0 & 0 & 20 & 10 & 0 & 0 & 40 & 30 \end{bmatrix}$$

doubly blocked =

$$\begin{bmatrix} \begin{matrix} 30 & 0 & 0 \\ 40 & 30 & 0 \\ 0 & 40 & 30 \\ 0 & 0 & 40 \end{matrix} & \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \\ \begin{matrix} 10 & 0 & 0 \\ 20 & 10 & 0 \\ 0 & 20 & 10 \\ 0 & 0 & 20 \end{matrix} & \begin{matrix} 30 & 0 & 0 \\ 40 & 30 & 0 \\ 0 & 40 & 30 \\ 0 & 0 & 40 \end{matrix} \\ \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} & \begin{matrix} 10 & 0 & 0 \\ 20 & 10 & 0 \\ 0 & 20 & 10 \\ 0 & 0 & 20 \end{matrix} \end{bmatrix}$$

⑥ Convert input matrix to a column vector.

$$I = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \Rightarrow \text{Vectorized } I = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

⑦ Multiply doubly blocked toeplitz matrix with vectorized ip signal.  
↳ This gives the convolution result.

⑧ Reshape the result to a matrix form.

$$\begin{array}{l} \text{Result vector} = \begin{bmatrix} 120 \\ 310 \\ 380 \\ 240 \\ 70 \\ 230 \\ 330 \\ 240 \\ 10 \\ 40 \\ 70 \\ 60 \end{bmatrix} \\ \Rightarrow \text{output} = \begin{bmatrix} 40 & 40 & 70 & 60 \\ 70 & 230 & 330 & 240 \\ 20 & 310 & 380 & 240 \end{bmatrix} \end{array}$$

2D (Cross) Correlation:

$$f[n, m] \otimes\otimes h[n, m] = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} f[n, m] * h[n+k, m+l] \rightarrow \text{Equivalent to convolution without the flip.}$$

↳ Compares the similarity of two sets of data.

↳ Measure of Relatedness of two signals.

# Template Matching, Gabor Filter and Edge Detection

## Gaussian Filter

$$(2D) G_{\sigma} = \frac{1}{2\pi\sigma^2} e^{-\left(\frac{x^2+y^2}{2\sigma^2}\right)}, \quad \sigma : \text{std. deviation}$$

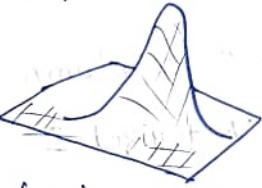
↳ determined extent of smoothing.

$$\left( \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) \right), \quad \begin{array}{l} \text{Infinite support, but discrete filters} \\ \text{use finite kernels.} \end{array}$$

$$\left( \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{y^2}{2\sigma^2}\right) \right), \quad \begin{array}{l} \text{Rule of thumb: set filter half-width} \\ \text{to about 3\sigma.} \end{array}$$

↓  
Product of 2 Gaussians

Gaussian:  $h_{\sigma}(x,y)$



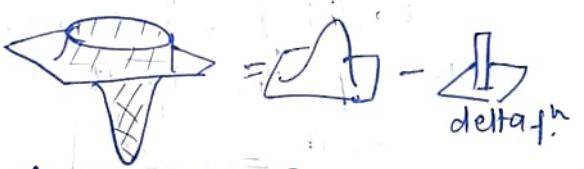
Derivative of Gaussian:  $\frac{\partial}{\partial x} h_{\sigma}(x,y)$



$$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

(Laplacian operator)

Laplacian of Gaussian:  $\nabla^2 h_{\sigma}(x,y)$



- Gaussian smoothing → for removing Gaussian noise.
- Reduces edge blurring
- Linear low-pass filter.
- Computationally efficient (implemented using small 1D filters)
- Rotationally symmetric.
- Degree of smoothing controlled by  $\sigma$ . ( $\sigma \uparrow$  for intensive smoothing.)
- To create soft shadow effects.

## Gabor Filter

- ↳ Product of complex sinusoid,  $c(x,y)$  [carrier]
- Gaussian function,  $g(x,y)$  [envelope].

$$c(x,y) = e^{-j[2\pi F(x\cos\theta + y\sin\theta) + \phi]}$$

$$g(x,y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-0.5\left(\frac{x_t}{\sigma_x^2} + \frac{y_t}{\sigma_y^2}\right)}$$

$$\begin{cases} x_t = (x\cos\theta + y\sin\theta)^2 \\ y_t = (-x\sin\theta + y\cos\theta)^2 \end{cases}$$

$$\text{gabor}(x,y) = c(x,y) * g(x,y) = e^{-\left(\frac{x^2+y^2}{2\sigma^2}\right)} \cos(2\pi\mu x) = \psi(x,y)$$

(2D)

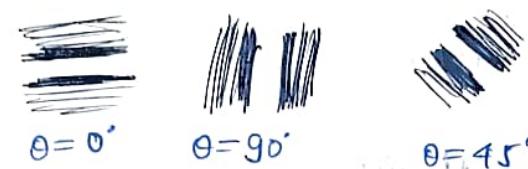
Canny:



Gaussian envelope:



Gabor kernel:



Median Filter → For removing salt and pepper noise

↳ Operates over a window by selecting the median intensity in the window.

Eg.

|    |    |    |  |  |  |  |  |  |
|----|----|----|--|--|--|--|--|--|
| 10 | 15 | 20 |  |  |  |  |  |  |
| 23 | 90 | 27 |  |  |  |  |  |  |
| 33 | 31 | 30 |  |  |  |  |  |  |

Sort

10 15 20 23 27 30 31 33 90

↑ Median value

↓ Replace

|    |    |    |  |  |  |  |  |  |
|----|----|----|--|--|--|--|--|--|
| 10 | 15 | 20 |  |  |  |  |  |  |
| 23 | 27 | 27 |  |  |  |  |  |  |
| 33 | 31 | 30 |  |  |  |  |  |  |

Template Matching

↳ Find part of one image (instance of pattern) that matched another.

↳ Use pattern as template.

↳ Convolution / correlation compares a template (filter) with each local image patch. → Dot product b/w filter & local image patch.

→ Normalized correlation → varies b/w -1 and 1.

↳ Attains 1 when filter & image region are identical.

## Edge Detection

- ↪ Identify sudden changes (discontinuities in an image)
- ↪ To extract information, recognize objects, recover geometry and viewpoint.

### Origin of edges:

- ① surface normal discontinuity
- ② depth discontinuity
- ③ surface color discontinuity
- ④ illumination discontinuity

## Image Gradients

### Types of discrete derivative in 1D:

Backward:  $\frac{df}{dx} = f(x) - f(x-1) = f'(x)$  filter:  $[0 \ 1 \ -1]$

Forward:  $\frac{df}{dx} = f(x) - f(x+1) = f'(x)$   $[-1 \ 1 \ 0]$

central:  $\frac{df}{dx} = \frac{f(x+1) - f(x-1)}{2} = f'(x)$   $[1 \ 0 \ -1]$

### Discrete Derivative in 2D: $f^n: f(x,y)$

Gradient vector,  $\nabla f = \begin{bmatrix} \frac{\partial f(x,y)}{\partial x} \\ \frac{\partial f(x,y)}{\partial y} \end{bmatrix} = \begin{bmatrix} f_x \\ f_y \end{bmatrix}$

Gradient magnitude,  $|\nabla f(x,y)| = \sqrt{f_x^2 + f_y^2}$

Gradient dirn,  $\theta = \tan^{-1} \left( \frac{\partial f}{\partial y} / \frac{\partial f}{\partial x} \right)$

3x3 Image gradient filter:  $\frac{1}{3} \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}, \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

### Characterizing Edges:

- ↪ Place of rapid change in image intensity function.

#### Effects of noise:

↪ Finite difference filters respond strongly to noise.

↪ Pixels look very different from their neighbours.

↪ Fix: smooth the image → Mean, Gaussian, Median filters.  
 $f(x)$ : 

$f'(x)$ :  → edge info lost

Solution : Smooth first

The diagram illustrates the convolution operation  $f * g$ . It shows three curves: a signal  $f$  (wavy line), a kernel  $g$  (narrow peak), and the result  $f * g$  (broad peak). A vertical dashed line marks the center of the kernel as it slides across the signal.

$$\frac{d}{dn} (f * g) = f * \frac{d}{dn} g$$

→ To find edges, look for peaks  
in  $\frac{d}{dx}(f * g)$ .

## Derivative of Gaussian Filter:

$$\text{2D Gaussian} \quad \times \star [0.1 \quad 0 \quad -1] =$$

→ Smoothed derivative removes noise, but blurs edges.

## Steps for edge detection:

## ① Select derivative filters

$$Sx = \begin{vmatrix} 1 & b & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{vmatrix}, \quad Sy = \begin{vmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ 1 & -2 & -1 \end{vmatrix}$$

② Convolve with image to compute derivatives.

$$\frac{\partial f}{\partial x} = s_x \otimes f, \quad \frac{\partial f}{\partial y} = s_y \otimes f$$

③ Form image gradient and compute its direction and amplitude.

Handling noise: Filter with a Gaussian to smooth, then take gradients.

Edge Adaptive Smoothing → ① Bilateral filter

## ② Non-local means filter

$I_p$ : value of image  $I$  at position  $p = (p_x, p_y)$ .

$F[I]$  = output of filter  $F$  applied to image  $I$ .

## Strategy for smoothing Images:

$$\text{Box Average: } BA[I]_p = \sum_{q \in S} \underbrace{B_\sigma(p-q)}_{\substack{\downarrow \\ \text{normalized}}} I_q \rightarrow \text{intensity at pixel } q.$$

→ Square box generates defects → Axis-aligned blocks

↳ Blocky results.

↳ Use an isotropic (circular) window

↳ Use a window with a smooth falloff.

↓ Gaussian Blur

$$GB[I]_P = \sum_{q \in S} G_{\sigma_s}(\|P-q\|) I_q$$

normalized gaussian fn.

↳ Smoother too much:  
edges are blurred.  
(bc of averaging across edge)

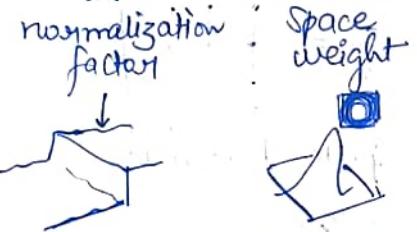
Bilateral Filter

↳ No averaging across edges.

↳ Kernel shape depends on image content.

$$BF[I]_P = \frac{1}{W_P} \sum G_{\sigma_s}(\|P-q\|) G_{\sigma_r}(|I_P - I_q|) I_q$$

normalization factor



(set, e.g., 2% of img. diagonal)



proportional to edge amplitude  
(set, e.g., mean or median of img. grads)

} Space  $\sigma_s$ : Spatial extent of kernel, size of considered neighborhood

} Range  $\sigma_r$ : "minm" amplitude of an edge.

} Independent of resolution and exposure.

Bilateral → Averages neighbor with similar intensities.

NL-Mean Filter

↳ Same goals: smooth within similar regions

↳ Generalize, extend 'similarity'

↳ Averages neighbors with similar neighborhoods.

→ For each and every pixel,  $P$ :

↳ define a small, simple fixed size neighborhood.

↳ define vector  $v_P$  → a list of neighboring pixel values.

$\sigma$  → size of window

↓ common setting:  
 $\alpha$  to image size

(e.g., 2% of  
image diagonal)

↳ independent of  
image resolution

Similar pixels  $p, q \Rightarrow$  small vector distance :  $\|\vec{v}_p - \vec{v}_q\|^2$

Dissimilar pixels  $p, q \Rightarrow$  large vector distance.

Filter with this: no spatial term.

$$NLMF[I]_p = \frac{1}{w_p} \sum_{q \in S} G_{\sigma_s}(\|\vec{v}_p - \vec{v}_q\|^2) I_q$$

vector distance to  
p gets weight for  
each pixel q.

Noisy source image

↪ Gaussian filter  $\rightarrow$  low noise, low detail

↪ Anisotropic Diffusion

↪ 'stairsteps'  $\sim$  piecewise constant

↪ Bilateral filter  $\rightarrow$  better, but similar stairsteps

↪ NL Means  $\rightarrow$  sharp, low noise, few artifacts.

Bilateral goals  $\rightarrow$  local smoothing within similar regions.

↪ Edge preserving smoothing.

↪ separate large structures & fine detail.

↪ Eliminate outliers.

↪ Filter within edges, not across them.

### Designing an Edge Detector

Criteria for "optimal" edge detector:

- ① Good detection  $\rightarrow$  Minimize probability of false positives (detecting spurious edges caused by noise), as well as false negatives (missing real edges).
- ② Good localization: The edges detected must be as close as possible to the true edges.
- ③ single response: The detector must return one point only for each true edge point, i.e., minimize the no. of local maxima around the true edge.

## Sobel Edge Detector

Sobel operator:

$$G_{x2} = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}, \quad G_{y2} = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} [1 \ 0 \ -1]$$

↑                      differentiation

↑                      Gaussian smoothing

Sobel operation:

Magnitude:  $G = \sqrt{G_x^2 + G_y^2}$

Dirn of gradient,  $\theta = \tan^{-1}\left(\frac{G_y}{G_x}\right)$

Sobel Filter problems:

- ↳ Poor localization (trigger response in multiple adjacent pixels)
- ↳ Threshold value favors certain dirns over others.
- ↳ can miss oblique edges more than horizontal or vertical edges.
- ↳ False negatives.

Unsharp Masking: Smoothing kernels used to sharpen images.

$$g_{\text{sharp}} = f + \gamma (f - h_{\text{blur}} * f)$$

Canny Edge Detector → Most widely used

- ① Suppress Noise (Gaussian smoothing & sobel grad. operator) [Detection]
- ② Compute gradient magnitude and direction (Sobel operator) [Localization]
- ③ Apply Non-Maximum Suppression (Localization)
  - ↳ Assumes minimal response.
- ④ Use hysteresis and connectivity analysis to detect edges.
  - ↳ Single edge point.

|   |   |
|---|---|
| $\text{Smooth} \rightarrow n_o * I$<br>$\text{Gradient} \rightarrow \nabla n_o * I$<br>$\text{Mag.} \rightarrow \  \nabla n_o * I \ $<br>$\text{Dirn.} \rightarrow \hat{n} = \frac{\nabla n_o * I}{\  \nabla n_o * I \ }$ | Compute Laplacian along<br>Grad. dirn $\hat{n}$ at each pixel:<br>$\frac{\partial^2 (n_o * I)}{\partial \hat{n}^2}$ |
|---|---|

## Non-Maximum Suppression:

- ↳ Edges occurs where gradient reaches a maxima.
- ↳ Suppress non-maxima gradient even if it passes threshold.
- ↳ Only 8 angle dims possible.
- ↳ Suppress all pixels in each dim which are not maxima.
- ↳ Do this in each marked pixel neighbourhood.

## Double threshold:

- ↳ High threshold to identify the strong pixels (intensity > high thresh.)
- ↳ Low threshold to identify non-relevant pixels (int. < low thresh.)
- ↳ All pixels having intensity in b/w are flagged as weak and Hysteresis mechanism helps identify the ones to be considered as strong and ones as non-relevant.

Hysteresis Threshold: Transform weak pixels into strong iff at least one of the pixels around the one being processed is a strong one.

- ↳ Declares an edge pixel if it is (weak) connected to an strong edge pixel ~~or~~ directly or via pixels b/w low & high.

## Hough Transform

- ↳ To find the location of lines in images.
- ↳ can detect lines, circles, and other structures only if their parametric eqn is known.
- ↳ Can give robust detection under noise & partial occlusion.

Primer to H.T.: Perform edge detection and a thresholding of some edge mag. image  $\rightarrow$  Have some pixels that may partially describe boundary of some objects.

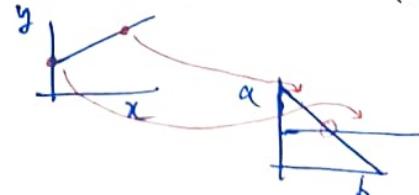
Find set of pixels that make a straight lines.

Point  $(x_i, y_i) \rightarrow$  Line passing through:  $y_i = ax_i + b$

$$\Rightarrow b = -ax_i + y_i \quad \hookrightarrow \text{satisfy eqn of some set of params } (a, b).$$

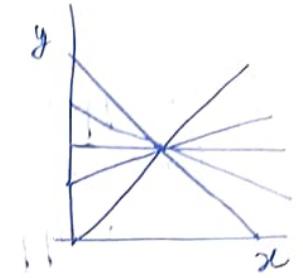
$x, y \rightarrow$  parameters  
 $a, b \rightarrow$  variables.

$\downarrow$   
A line in  $(a, b)$  space



2 points  $(x_1, y_1), (x_2, y_2) \rightarrow$  define a line in  $(x, y)$  plane.

↳ Given 2 different lines in  $(a, b)$  space.  
↓  
intersect at  $(a', b')$ .



$\Rightarrow$  All pts. on line defined by these two points in  $(a, b)$  space will parameterize lines that intersect in  $(a', b')$  in  $(a, b)$  space.

Algorithm for H.T.:

- ① Quantize the parameter space  $(a, b)$  by dividing it into cells.
- ② Count the no. of times a line intersect a given cell.  $\Rightarrow$  Accumulator cells.
- ③ Cells receiving more than a certain no. of counts (called votes) are assumed to correspond to lines in  $(x, y)$  space.

### Integral Image (Summed Area Table)

↳ If an image is to be repeatedly convolved with different box filters (and especially of different sizes at diff. locations), we can precompute the summed area table, which is just the running sum of all pixel values from the origin.

$$S(i, j) = \sum_{k=0}^i \sum_{l=0}^j f(k, l)$$

Recursive (Raster-scan) algorithm:

$$S(i, j) = S(i-1, j) + S(i, j-1) - S(i-1, j-1) + f(i, j)$$

Eg.

|   |   |   |   |
|---|---|---|---|
| 3 | 2 | 7 | 2 |
| 1 | 5 | 1 | 3 |
| 5 | 1 | 3 | 5 |

T/P Image

|   |    |    |    |
|---|----|----|----|
| 3 | 5  | 12 | 14 |
| 4 | 11 | 19 | 24 |
| 9 | 17 | 28 | 38 |

Summed Area Table

$$28 = 3 + 19 + 17 - 11$$

Summed Area:

$$S(i_0 \dots i_1, j_0 \dots j_1) = S(i_1, j_1) - S(i_1, j_0 - 1) - S(i_0 - 1, j_1) + S(i_0 - 1, j_0 - 1)$$

$$\begin{aligned} S &= 5 + 1 + 3 = 28 + 3 - 12 - 9 \\ &= 10 \end{aligned}$$

# Frequency Domain Analysis of Images

$$\text{Fourier transform: } F(f) = \int_{-\infty}^{\infty} f(x) e^{-j2\pi fx} dx$$

$$\text{Inverse FT : } f(x) = \int_{-\infty}^{\infty} F(u) e^{j2\pi ux} du$$

2D FT:

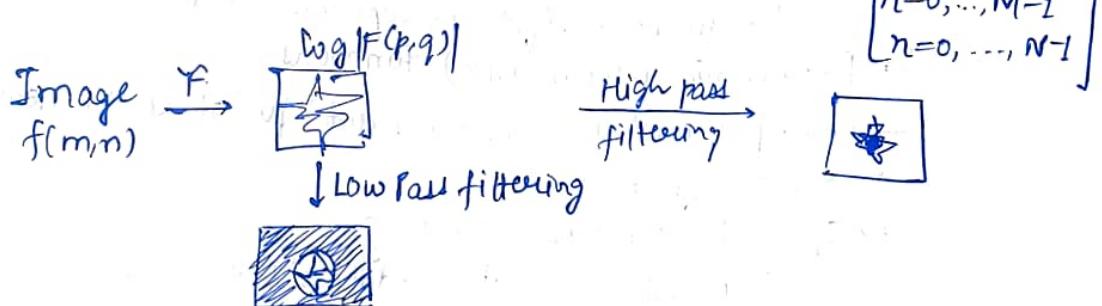
$$\text{FT: } F(u, v) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) e^{-j2\pi(ux+vy)} dx dy$$

$$\text{IFT: } f(x,y) = \iint_{-\infty}^{\infty} F(u,v) e^{j2\pi(xu+yv)} du dv$$

$$\text{Discrete FT: } F[p, q] = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} f[m, n] e^{-j2\pi pm/M} e^{-j2\pi qn/N}$$

$p, q$ : frequencies along  $m, n$   $\begin{cases} p=0, \dots, M-1 \\ q=0, \dots, N-1 \end{cases}$

$$\text{IDFT: } f[m,n] = \frac{1}{MN} \sum_{p=0}^{M-1} \sum_{q=0}^{N-1} F[p,q] e^{j2\pi p m / M} e^{j2\pi q n / N}$$



→ FT has peaks at spatial frequencies of repeated structure.

# Deep Learning

Artificial Intelligence → scientific field concerned with the development of algorithms that allow computers to learn without being explicitly programmed.

Machine Learning → Branch of AI which focuses on methods that learn data and make predictions on unseen data.

Labelled data → ML algorithm

Training  
Prediction

Labeled data → Learned model → Prediction

## ML Types

① Supervised: Learning with labeled data.

② Unsupervised: Discover patterns in unlabeled data.

③ Reinforcement learning: Learn to act on feedback/reward.

→ Supervised learning categories and techniques:

↳ Numerical classifier functions  
↳ Linear classifier, perceptron, logistic regression, SVM, NN.

↳ Parametric (probabilistic) functions

↳ Non-parametric (instance-based) functions → k-nearest neighbors.

↳ Symbolic functions.

↳ Aggregation (ensemble) learning

↳ Bagging, boosting (Adaboost), random forest.

→ Unsupervised learning categories and techniques:

↳ Clustering

↳ K-Means, Mean-shift, Spectral.

↳ Density estimation.

↳ Dimensionality reduction.

↳ Principal component analysis (PCA), factor analysis.

Nearest Neighbor classifier: For each data point, assign the class label of the nearest training data point.

↳ Distance function to find nearest neighbor.

↳ Does not require learning a set of weights.

↳ Remember all training data and store.

## K-Nearest Neighbors Classifier:

- ↳ Consider multiple neighbouring data points  $(k)$  to classify a test data point.
- ↳ Class of test eg. is obtained by voting.

## Linear classifier:

- ↳ Find a function  $f$  of the inputs  $x_i$  that separates the classes.

$$f(x_i, w, b) = w x_i + b.$$

- ↳ Use pairs of inputs and labels to find weight matrix  $w$  and bias vector  $b$  (parameters of  $f$ )
- ↳ Several methods to find optimal set of parameters: eg., perceptron.
- ↳ Decision boundary is linear.

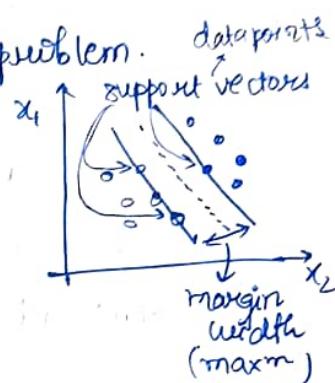
## Support Vector Machines (SVM): Solves an optimization problem.

- ↳ Identify a decision boundary that correctly classifies the examples.
- ↳ Then, increase the geometric margin b/w the boundary & all examples.

Find  $w$  and  $b$  by solving:

$$\min \frac{1}{2} \|w\|^2$$

such that  $y_i(w \cdot x_i + b) \geq 1 \forall x_i$ .



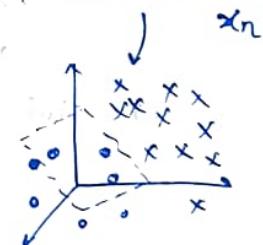
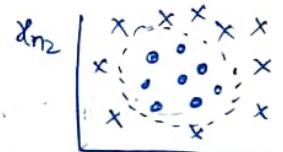
## Non-linear Techniques: Features $z_i$ are obtained as non-linear functions of the i/p's $x_i$ .

- ↳ Non-linear decision boundaries.

## Non-linear SVM:

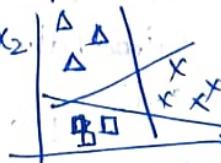
- ↳ The original i/p space is mapped to a higher-dimensional space where the training set is linearly separable.

- ↳ Define a non-linear kernel function to calculate a non-linear decision boundary in the original feature space.



## Multi-class classifications:

- ↳ 3 or more classes.



## Computer Vision Tasks: Classification, localization, object detection, instance segmentation.

## Non-Free-Lunch Theorem: No single classifier works the best for all possible problems.

- ↳ Need to make assumptions to generalize.

Deep Learning: Uses multiple layers for learning data representation.

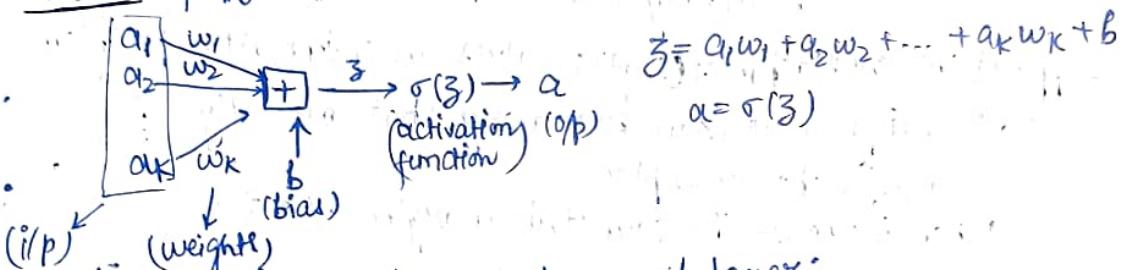
ML: I/p → Feature extraction → classification → output

DL: I/p → Feature extraction → classification → Output

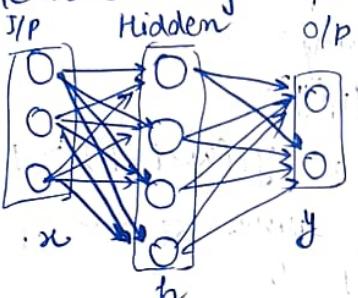
→ Neural networks with at least one hidden layer are universal approximators.

↳ Use non-linear mapping of the inputs  $x$  to the outputs  $f(x)$  to compute complex decision boundaries.

Neuron:  $f: \mathbb{R}^k \rightarrow \mathbb{R}$



→ NN with one hidden layer and one o/p layer:



$$\text{hidden layer: } h = \sigma(w_1 x + b_1)$$
$$\text{o/p layer: } y = \sigma(w_2 h + b_2)$$

4+2=6 neurons (not counting bias)

$$\{(8 \times 4) + (4 \times 2)\} = 20 \text{ weights}$$

$$4+2=6 \text{ biases}$$

↳ 26 learnable parameters

→ Fully connected (dense) layers

↳ Multilayer perceptron (MLP)

↳ Each neuron is connected to all neurons in succeeding layer.

→ Matrix operation:

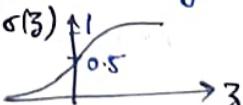
$$\sigma(W[x] + b) = a$$

Softmax layer: In multi-class classification

↳ Employs softmax activation function → Range [0, 1].

↳ Sigmoid activations → for binary classification

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$$y_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

↳  $z$  inputted to softmax layer: logits  
↳ probability:  $0 < y_i < 1$

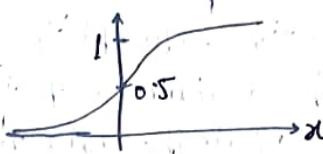
## Activation Functions:

- ↳ Non-linear activations are needed to learn complex (non-linear) data representations.
- ↳ Otherwise, NNs would be just a linear function ( $w_1 w_2 x = w_2 x$ ).

Eg. ① sigmoid

$$f(x) = \frac{1}{1+e^{-x}}$$

$[0, 1]$



↳ Less common:  $\mathbb{R}^n \rightarrow [0, 1]$

② Tanh

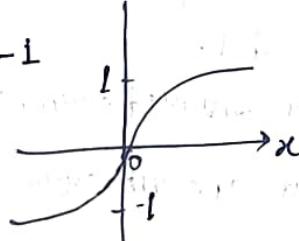
$$\tanh(x) = \frac{2}{1+e^{-2x}} - 1$$

↳ Saturates

↳ zero-centered

$$\tanh(x) = 2 \cdot \sigma(2x) - 1$$

↳ sigmoid



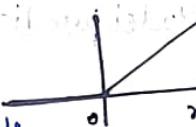
$\mathbb{R}^n \rightarrow [-1, 1]$

③ ReLU (Rectified Linear Unit)

$$f(x) = \max(0, x)$$

$$f(x) = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$

↳ Used in most modern deep NNs.



$\mathbb{R}^n \rightarrow \mathbb{R}^n$

↳ Fast to compute

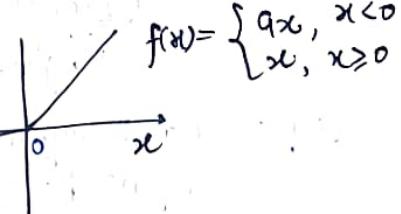
↳ Accelerates the convergence of gradient descent

↳ Prevents gradient vanishing problem

④ Leaky ReLU

↳ Has a small -ve slope  
(e.g.,  $\alpha = 0.01$ )

↳ Resolves dying ReLU problem.



$$f(x) = \begin{cases} ax, & x < 0 \\ x, & x \geq 0 \end{cases}$$

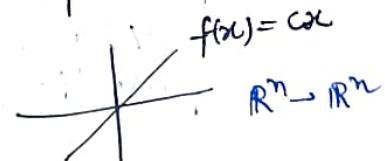
$\mathbb{R}^n \rightarrow \mathbb{R}^n$

⑤ Linear function

↳ O/p  $\propto$  I/p signal

↳  $c \approx 1$ : identity activation fn

↳ Used in regression problems.



$$f(x) = cx$$

$\mathbb{R}^n \rightarrow \mathbb{R}^n$

## Training NNs

Network parameter,  $\theta = \{w^1, b^1; w^2, b^2, \dots, w^L, b^L\} \rightarrow$  weight matrices  
 bias vectors

→ Train the model to learn a set of parameters  $\theta$  that are optimal.

① Data preprocessing → Helps convergence during training.

↳ Mean subtraction → to obtain zero-centered data.

↳ Normalization → divide each feature by its std. deviation.  
 ↳ scale the data within range  $[0, 1]$  or  $[-1, 1]$ .

② Define a Loss / objective / cost function,  $L(\theta)$

↳ calculate difference (error) b/w model prediction and true label.

$$\text{Eq. } \left\{ \begin{array}{l} \text{(i) Mean squared error: } L(\theta) = \frac{1}{n} \sum_{i=1}^n [y^{(i)} - \hat{y}^{(i)}]^2 \\ \text{(ii) Mean absolute error: } L(\theta) = \frac{1}{n} \sum_{i=1}^n |y^{(i)} - \hat{y}^{(i)}| \end{array} \right.$$

$$\text{(iii) Cross-entropy: } L(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K [y_k^{(i)} \log \hat{y}_k^{(i)} + (1-y_k^{(i)}) \log (1-\hat{y}_k^{(i)})]$$

for classification tasks

Ground truth class labels  $\rightarrow y_i$   
 Model predicted class labels  $\rightarrow \hat{y}_i$

## Gradient Descent Algorithm

↳ To optimize loss function,  $L(\theta)$ .

↳ Applies iterative refinement of the n/w parameters  $\theta$ .

↳ Use the opposite dirn. of the gradient of the loss w.r.t. NN param.

(i.e.,  $\nabla L(\theta) = \partial L / \partial \theta$ ) for updating  $\theta$ .

↳ give the dirn. of fastest increase of  $L(\theta)$  when  $\theta$  are updated.

Steps:

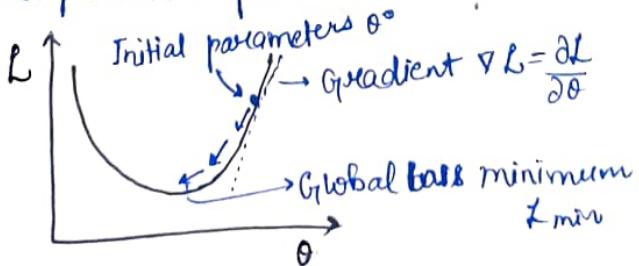
- ① Randomly initialize the model parameters,  $\theta^0$ .

② compute the gradient of the loss function at initial parameters  $\theta^0$ :  $\nabla L(\theta^0)$ .

③ Update the parameters as:  $\theta^{\text{new}} = \theta^0 - (\alpha \nabla L(\theta^0))$

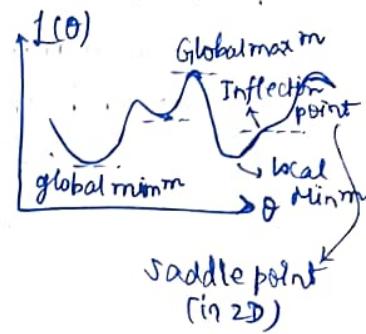
$[\alpha: \text{learning rate}]$

④ Go to step 2 and repeat / until a termination criterion is reached.



$$\nabla L(\theta^0) = \begin{bmatrix} \frac{\partial L(\theta^0)}{\partial \theta_1} \\ \frac{\partial L(\theta^0)}{\partial \theta_2} \\ \vdots \end{bmatrix}$$

- GD algo. stops when a local minimum of the loss surface is reached.
- Does not guarantee reaching global minimum.
- However, empirical evidence suggests that GD works well for NNs.
- Random initialization in NNs results in different initial parameters  $\theta$  every time the NN is trained.
  - ↳ GD may reach different minima at every run.
  - ↳ Produce different predicted o/p.



Forward propagation: (forward pass) Passing the i/p  $x$  through the hidden layers to obtain model o/p (prediction)  $y \rightarrow$  calculate  $L$ .

Backpropagation: (backward propagation)

- ↳ Traverses the n/w in reverse order, from o/p  $y$  backward towards the i/p  $x$  to calculate the gradients of the loss  $\nabla L(\theta)$ .
- ↳ chain rule is used to calculate the partial derivatives of the loss fn w.r.t.  $\theta$  in different layers in the n/w.
- ↳ Each update of model parameter  $\theta$  during training takes one forward and one backward pass (e.g., of a batch of i/p's).
- ↳ Automatic calculation of gradients (automatic differentiation) is available in all current deep learning libraries.

Mini-Batch G.D.:

- ↳ compute loss  $L(\theta)$  on a mini-batch of images, update the parameters  $\theta$  and repeat until all images are used.
- ↳ Much faster training.
- ↳ Works b/c the gradient from a mini-batch is a good approximation of the gradient from the entire training set.

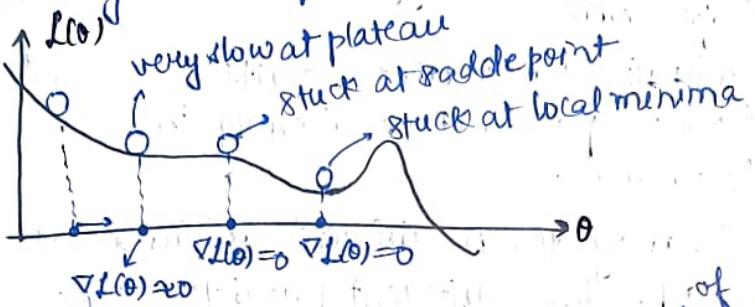
Stochastic Gradient Descent:

- ↳ Use mini-batches that consists of a single input example ( $I_{image}$ ), e.g.
- ↳ Fast, but causes significant fluctuations in the loss function.

## Problems with G.D.

↳ Local minima

↳ can be very slow at plateaus and can get stuck at saddle points.



G.D with momentum: Use momentum ~~for~~ gradient for optimization.

Weighted average of past gradients

Movement = -ve of gradient + Momentum

$$\theta^t = \theta^{t-1} - v_t, \quad v_t: \text{momentum} \rightarrow \text{accumulated grad. from past several steps}$$

$$\text{where } v_t = \beta v^{t-1} + \alpha \nabla L(\theta^{t-1})$$

$$\therefore \theta^t = \theta^{t-1} - \alpha \nabla L(\theta^{t-1}) - \beta v^{t-1}$$

prev. iteration

$\beta$ : coeff. of momentum  
↳ typically = 0.9

Nesterov Accelerated Momentum:

$$\theta^t = \theta^{t-1} - v_t$$

$$\text{where } v_t = \beta v^{t-1} + \alpha \nabla L(\theta^{t-1} + \beta v^{t-1})$$

↳ allows to predict position of param. in the next step.

Adam (Adaptive Moment Estimation)

weighted avg.  $\theta^t = \theta^{t-1} - \alpha \frac{\hat{v}^t}{\sqrt{\hat{v}^t + \epsilon}}$ , where  $\hat{v}^t = \frac{v^t}{1-\beta_1}, \hat{s}^t = \frac{s^t}{1-\beta_2}$

offset and weighted avg. of past squared gradients.

Typically,  $\beta_1 = 0.9$   
 $\beta_2 = 0.999$   
 $\epsilon = 10^{-8}$

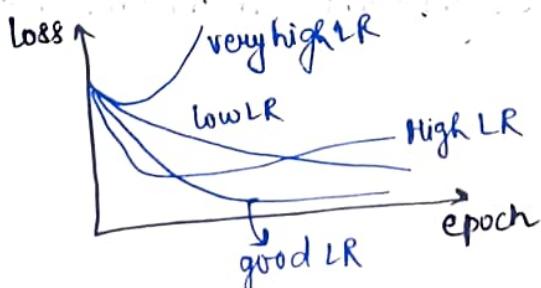
$$v^t = \beta_1 v^{t-1} + (1-\beta_1) \nabla L(\theta^{t-1})$$

$$s^t = \beta_2 s^{t-1} + (1-\beta_2) (\nabla L(\theta^{t-1}))^2$$

Learning Rate: Step size

High LR  $\Rightarrow$  Loss  $\uparrow$  or plateaus too quickly

Low LR  $\Rightarrow$  Loss  $\downarrow$  too slowly (takes many epochs)



## Learning Rate Scheduling: Change LR during training

- ↳ Annealing: Reduces LR over time (LR decay / static decay).
- ↳ Approaches:
  - ↓ LR by some factor every few epochs.
  - Exponential or cosine decay over time.
  - ↓ LR by a constant (e.g., by half) whenever validation loss stops improving.
- ↳ Warmup: Gradually ↑ LR initially, afterward let it cool down until the end of training.

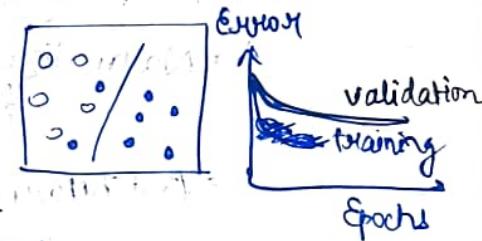
## Vanishing Gradient Problem

- ↳ During training, gradients can become either very small (vanishing gradient) or very large (exploding gradient).
- ↳ Very small or very large update of parameters.
- ↳ Solutions: change LR, ReLU activations, regularization, LSTM units in RNNs.

## Generalization:

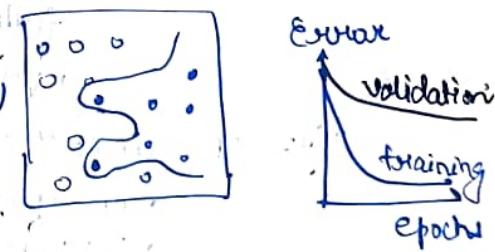
### ① Underfitting

- ↳ Model too simple to represent all relevant class characteristics.
- ↳ Model with too few parameters.
- ↳ Produces high error on training set and high error on validation set.



### ② Overfitting:

- ↳ Model too complex and fits irrelevant characteristics (noise) in the data.
- ↳ Model with too many params.
- ↳ Produces low error on training error and high error on the validation set.



## Regularization: Weight Decay

↳ weight decay: Penalizes large weights added to the loss function

$$L_{\text{reg}}(\theta) = L(\theta) + \lambda \sum_k \theta_k^2$$

$\lambda$ : weight decay coeff.  
determines how dominant regularizer is during gradient computation.

$\ell_1$  weight decay: Less common  $\rightarrow$  performs worse than  $\ell_2$  decay.

$$L_{\text{reg}}(\theta) = L(\theta) + \lambda \sum_k |\theta_k|$$

$\rightarrow$  combine  $\ell_1$  and  $\ell_2$  regularization (Elastic net regularization).

$$L_{\text{reg}}(\theta) = L(\theta) + \lambda_1 \sum_k |\theta_k| + \lambda_2 \sum_k \theta_k^2$$

Regularization: Dropout

$\hookrightarrow$  Randomly drop units (along with their connections) during training. [Often 20% - 50% units dropped]

$\hookrightarrow$  A kind of ensemble learning.

Regularization: Early Stopping

$\hookrightarrow$  During model training, use of validation set.

$\hookrightarrow$  Stop when validation accuracy (or loss) has not improved after n epochs. [ $n \rightarrow$  patience parameter].

Batch Normalization  $\rightarrow$  similar to data preprocessing step.

$\hookrightarrow$  Normalize data  $x$  to zero mean and unit variance.

$$\hat{x} = \frac{x - \mu}{\sigma}$$

$\hookrightarrow$  BatchNorm layers alleviate the problems of proper initialisation of parameters and hyper-parameters.

$\hookrightarrow$  Faster convergence training, allow  $\uparrow$  LR.

$\hookrightarrow$  suited immediately after convolutional layers or fully-connected layers, and before activation layers.

$\hookrightarrow$  very common with CNNs.

Hyper-Parameter Tuning: Setting hyper parameters.

$\hookrightarrow$  No. of layers, no. of neurons per layer

$\hookrightarrow$  Initial LR, LR decay schedule (eg., decay constant).

$\hookrightarrow$  Optimizer type

$\hookrightarrow$  Regularization parameters

$\hookrightarrow$  Batch size

$\hookrightarrow$  Activation functions

$\hookrightarrow$  Loss function

$\hookrightarrow$  Grid search: check all values in a range with step value.

$\hookrightarrow$  Random search: Randomly sample values for parameters.

$\hookrightarrow$  Bayesian hyper-parameter optimization.

K-Fold Cross-Validation → Hyper-parameter tuning when size of training data is small.  
↳ Better and less noisy estimate of model performance by averaging the results across several folds.

Eg. 5-fold CV.

↳ Split train data into 5 equal folds.

① First use 2-5 for training and fold 1 for validation.

② Then use 1-5 for training and fold 2 for validation, then fold 3, 4, 5.

③ Repeat by using fold 2 for validation, then fold 3, 4, 5.

④ Average the results over 5 rounds.

⑤ Determine best hyper-param and evaluate model on test data.

## Ensemble Learning

↳ Training multiple classifiers separately and combining their predictions. → Outperforms individual classifiers.

↳ Bagging (bootstrap aggregating):

↳ Randomly draw subsets from training set (bootstrap samples).

↳ Train separate classifiers on each subset of training set.

↳ Perform classification based on average vote of all classifiers.

↳ Boosting:

↳ Train a classifier and apply weights on the training set (apply higher weights on misclassified examples).

↳ Train new classifier, reweight training set a/c to prediction error.

↳ Repeat.

↳ Perform classification based on weighted vote of classifiers.

## Deep vs Shallow Networks

↳ Deeper nets perform better than shallow nets, but only upto some limit (after that performance plateaus).

## Convolutional Neural Networks (CNNs)

- ↳ Primarily designed for image data.
- ↳ Use a convolutional operator for extracting data features.
  - ↳ Allow parameter sharing.
  - ↳ Efficient to train.
    - ↳ Have less parameters than fully-connected NNs.
- ↳ Robust to spatial translations of objects in images.
  - ↳ A convolutional filter slides (convolves) across the image.
  - ↳ Hidden units in a layer are only connected to a small region of the layer before it (called receptive field).
    - capture useful features.
  - ↳ Feature map: Depth of each feature map corresponds to the no. of convolutional filters used at each layer.
- ↳ Max pooling: Reports the max o/p within a rectangular region.
- ↳ Average pooling: Reports the avg o/p of a rectangular neighborhood.
- Pooling layers reduce the spatial size of feature maps.
  - ↳ Reduce no. of params., prevent overfitting.

## RANSAC and Feature Detectors

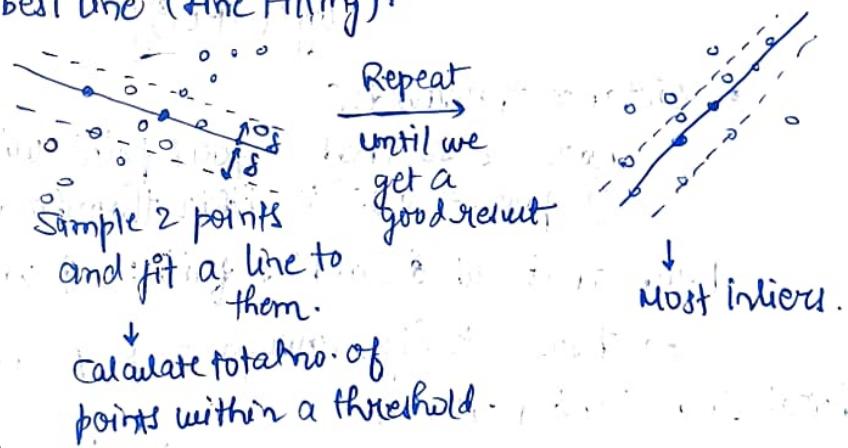
### RANSAC (RANDOM Sample Consensus)

- ↳ Robust model fitting algorithm that can estimate a model from data containing many outliers.
- ↳ For edge detection, line fitting, homography estimation, 3D reconstruction.

#### RANSAC Loop:

- ① Randomly select a seed group of points on which to base transformation estimate.
- ② Compute transformation / Fit a model using only sampled pts.
- ③ Find inliers.
- ④ If the no. of inliers is sufficiently large, re-compute least squares estimate of transformations on all of the inliers.  
→ Keep the model/transformations with the largest no. of inliers.

Eg. Estimate best line (Line Fitting):



### Feature-Based Alignment:

- ① Extract features.
- ② Compute putative matches.
- ③ Loop:
  - Hypothesize transformation  $T$ .
  - Verify transformation (search for other matches consistent with  $T$ ).

## Choosing parameters:

- No. of sample points :  $s$  (2 for fitting a line).
- Outlier ratio :  $e = \frac{\text{no. of outliers}}{\text{no. of datapoints}}$
- No. of trials :  $T$

Choose  $T$  so that, with probability  $p$ , at least one random sample is set free from outliers.

$p(\text{fail once}) = \text{don't select any outlier}$

$$\Rightarrow 1-p = [1-(1-e)^s]^T \quad | \quad \begin{array}{l} (1-e)^s: \text{In one trial, prob.} \\ \text{prob. that a sample} \\ \text{contains at least one} \\ \text{outlier (ie, is bad)} \end{array}$$

$p(\text{fail } T \text{ times}) = \text{select at least one outlier in all } T \text{ trials}$

$$1-p = [(1-(1-e)^s)]^T$$

$$\Rightarrow \log(1-p) = T \log(1-(1-e)^s)$$

$$\Rightarrow T = \frac{\log(1-p)}{\log(1-(1-e)^s)}$$

- Distance threshold:  $\delta$

choose  $\delta$  so that a good point with noise is likely (eg, prob. = 0.95) is within threshold.

Prob. that all  $K$  samples fail :  $(1-w^\delta)^K$  ||  $K$ : samples chosen

$\Rightarrow$  choose  $K$  high enough to keep this below desired failure rate.

## Local Invariant Features

$\hookrightarrow$  Global representations have major limitations  $\Rightarrow$  Match local regions.

$\hookrightarrow$  Increased robustness to occlusion, articulation, intra-category variations.

- steps:
- ① Find a set of distinctive key-points
  - ② Define a region around each keypoint
  - ③ Extract and normalize the region content
  - ④ Compute a local descriptor from normalized region
  - ⑤ Match local descriptors

$\hookrightarrow$  Region extraction needs to be repeatable and accurate.

$\hookrightarrow$  Geometric & photometric invariant.

## Harris Corner Detection

↳ Identifies corner in an image by measuring how the intensity changes in all directions around a pixel.



$$\sum I_x^2 \rightarrow \text{large}$$

$$\sum I_y^2 \rightarrow \text{large}$$

(corner)

$$\sum I_x^2 \rightarrow \text{small}$$

$$\sum I_y^2 \rightarrow \text{large}$$

(Edge)

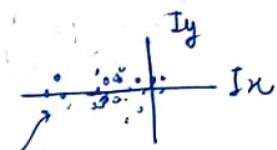
$$\sum I_x^2 \rightarrow \text{small}$$

$$\sum I_y^2 \rightarrow \text{small}$$

(Nothing)

→ flat

↓  
Small movement  
in any dirn causes  
large change in intensity



Finding corners:

- ① compute image gradients over small region.

② subtract mean from each image gradient:

③ compute covariance matrix.

④ compute eigenvectors and eigenvalues.

⑤ use threshold on eigenvalues to detect corners.

$$I_x = \frac{\partial I}{\partial x}, I_y = \frac{\partial I}{\partial y}$$

centred data

Covariance matrix:

$$\begin{bmatrix} \sum_{p \in P} I_x I_y & \sum_{p \in P} I_x^2 \\ \sum_{p \in P} I_y I_x & \sum_{p \in P} I_y^2 \end{bmatrix}$$

Intensity change for change in intensity  $E(u, v)$  when window is shifted by  $(u, v)$  is

$$E(u, v) = \sum_{x, y} w(x, y) [I(x+u, y+v)] - [I(x, y)]^2$$

First-order Taylor expansion:

$$I(x+u, y+v) \approx I(x+y) + u I_x + v I_y$$

$$\therefore E(u, v) \approx \sum_{x, y} w(x, y) [u I_x + v I_y]^2$$

$$= [u, v] \sum_{x, y} w(x, y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

# Compute eigenvalues & eigenvectors

$$M\mathbf{e} = \lambda \mathbf{e}$$

$$(M - \lambda I)\mathbf{e} = 0$$

$M \rightarrow$  autocorrelation matrix

$$\hookrightarrow ① \det(M - \lambda I) = 0 \Rightarrow \lambda_1, \lambda_2$$

$$② \text{For } \lambda_1, \lambda_2, \text{ solve } (M - \lambda I)\mathbf{e} = 0$$

$\hookrightarrow 2 \times 2$  Harris matrix

## Corner Response Function

$M \rightarrow$  eigenvalues:  $\lambda_1, \lambda_2$

$$R = \det(M) - K(\text{trace}(M))^2$$

$$= (\lambda_1 \lambda_2) - K(\lambda_1 + \lambda_2)^2$$

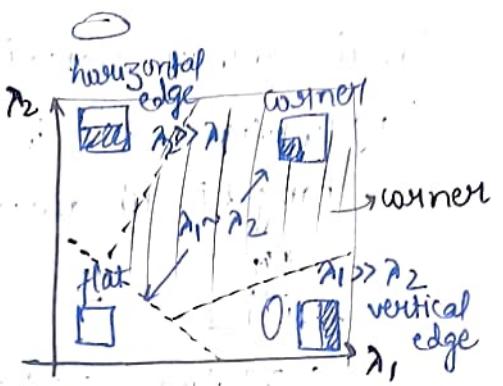
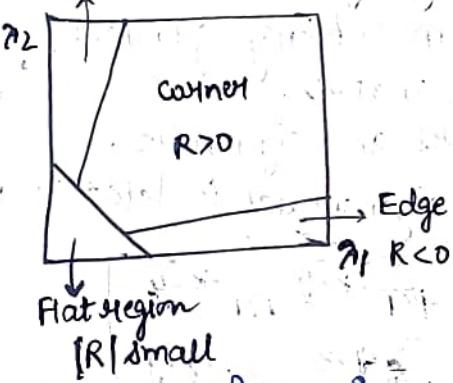
$M \rightarrow$  symmetric

$$M = R^{-1} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} R$$

Ellipse fm:

$$[U \ V] M \begin{bmatrix} U \\ V \end{bmatrix} = \text{const.}$$

Edge ( $R < 0$ )



# Kanade & Tomasi:

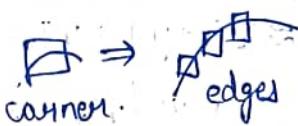
$$R = \min(\lambda_1, \lambda_2)$$

# Nobel:

$$R = \frac{\det(M)}{\text{trace}(M) + \epsilon}$$

## Harris Detector: Properties

- ↳ Translation invariance
- ↳ Rotation invariance
- ↳ Partial invariance to affine intensity change
- ↳ Not scale invariant



## Histograms of Oriented Gradients (HoG) and

## Bag of Visual Words (BoVW) features

### SIFT (Scale-Invariant Feature Transform) → Advantages:

- Locality: features are local, robust to occlusion and clutter.
- Distinctiveness: individual features can be matched to a large database of objects.
- Quantity: many features can be generated for even small objects.
- Efficiency: clear for real-time performance.
- Extensibility: To different feature types.

### SIFT Algorithm:

- ① scale-space peak selection: potential location of finding features.
- ② Keypoint localization: Accurately locating feature keypoints.
- ③ Orientation assignment: Assigning orientation to keypoints.
- ④ Keypoint descriptor: Describing keypoints as a high-dimensional vector.
- ⑤ Keypoint matching.

### HoG:

- ① Compute a HoG.
- ② (Optional) global image localization.
- ③ Computing gradient image in  $x$  and  $y$ .
- ④ Computing gradient histograms.
- ⑤ Normalizing across blocks.
  - ↳ Focuses on structure or shape of an object.
  - ↳ Orientations are calculated in localized patches portions.
- ⑥ Flattening into a feature vector.

### Step: ① Preprocess the data

Width to height ratio = 1:2 (preferably 64x128).

↳ For dividing images into 8x8 and 16x16 patches.

- ② calculate gradients (dim  $x$  and  $y$ )
- ③ calculate magnitude and orientation.
- ④ calculate histogram of gradient in 8x8 cells.
- ⑤ Normalize gradients in 16x16 cell (36x1)  $\rightarrow$  Reduce lighting variation
- ⑥ Features for complete image.
  - ↳ Combine 16x16 blocks to get features for the final image.

## Bag of Words

- visual words = independent features
- words from the dictionary.
- Represent images based on histogram of word occurrences.
- Image comparisons are performed based on such word histograms.

### ① Feature detection and representation

sliding window → visual words (keypoints)  
 regular grid → k-Means clustered pieces of image  
 interest point detector

- Image to histogram
- Dictionary → list of words → learned from data
- Extract feature descriptors from a training dataset.
  - points in high-dimensional space
  - ↳ cluster similar descriptors (k-Means)

Task: Find similar looking images

Input: Database of images, dictionary, query image(s).

Output: N most similar database images  $\xrightarrow{\text{to}}$  query images.  
 → words that occur in every image do not help a lot for comparisons.

TF-IDF Reweighting: of word in a document, of word across several documents

TF-IDF = term frequency - inverse document frequency

$$tid = \frac{n_{id} \log N}{n_d}, \quad tid: \text{histogram bin of word } i \text{ for image } d$$

term frequency IDF  $n_{id}$ : occurrences of word  $i$  in image  $d$

→ Relevant words get higher weights, others are weighted down to 0. (those occurring in every image).  
 $n_d$ : no. of word occurrences in  $d$ .  
 $n_i$ : no. of images that contain word  $i$ .  
 $N$ : no. of images

Comparing two Histograms:

$$\text{cosine similarity: } \text{cosim}(x, y) = \cos \theta = \frac{x^T y}{\|x\| \|y\|} \quad [0, 1]$$

$$\text{cosine distance: } d_{cos}(x, y) = 1 - \text{cosim}(x, y) = 1 - \frac{x^T y}{\|x\| \|y\|} \quad \Rightarrow \text{form cost matrix}$$

$$\text{Euclidean distance: } \|x - y\|^2 = 2 d_{cos}(x, y) \quad [\text{for } \|x\| = \|y\|]$$

# Video Processing

Optical Flow: Given two consecutive image frames, estimate the motion of each pixel.

↪ Assumption: Brightness constancy, small motion.

Color/Brightness Constancy: Allowed for pixel-to-pixel comparison (not image features)

$$I(x(t), y(t), t) = C$$

$$\frac{dI}{dt} = \frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{dy}{dt} + \frac{\partial I}{\partial t} = 0 \quad \dots \text{Brightness constancy eqn.}$$

Derivation:

$$I(x+t\delta t, y+v\delta t, t+\delta t) = I(x, y, t)$$

Assuming small motion, Taylor series expansion:

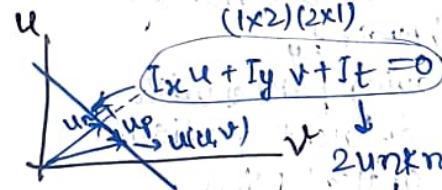
$$I(x, y, t) = I(x, y, t) + \frac{\partial I}{\partial x} \delta x + \frac{\partial I}{\partial y} \delta y + \frac{\partial I}{\partial t} \delta t$$

Divide by  $\delta t$ . Take limit  $\delta t \rightarrow 0$ :

$$\Rightarrow \frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{dy}{dt} + \frac{\partial I}{\partial t} = 0$$

$$\text{or, } I_x u + I_y v + I_t = 0$$

$$\text{or, } \nabla I^T \cdot \begin{pmatrix} u \\ v \end{pmatrix} + I_t = 0$$



$$u = u_n + u_p$$

normal flow      parallel flow

2 unknowns  
(unknown)

2 equations  
needed

compute by Sobel filter,  
DOG filter, forward  
difference

$I_x, I_y$ : image gradients  
(spatial derivatives)

$u, v$ : flow velocities

$I_t$ : temporal gradient

optical flow  
flow

compute by frame  
difference

solve for this  
constraint  
→ 2nd lies on a line.

cannot be found uniquely with  
a single constraint.

$$\hat{u}_n = \frac{(I_x, I_y)}{\sqrt{I_x^2 + I_y^2}}, \quad |u_n| = \frac{|I_t|}{\sqrt{I_x^2 + I_y^2}} \Rightarrow \hat{u}_n = \frac{|I_t|}{\sqrt{I_x^2 + I_y^2}} (I_x, I_y)$$

$u_p \rightarrow$  cannot be determined

Small Motion: Linearization of brightness constancy constraint

## Horn-Schunck Optical Flow

- Brightness constancy
- small motion
- smooth flow  
(flow can vary from pixel to pixel)
- Global method (dense)

## Lucas-Kanade Optical Flow

- Method of differences
- constant flow  
(flow is constant for all pixels)
- Local method  
(sparse)

## Constant Flow: Lucas-Kanade Optical Flow

$$I_x u + I_y v + I_t = 0$$

Given 25 equations.

→ Assume that the surrounding path (say 5x5) has constant flow.

→ Assumption: Flow is locally smooth.

Neighbouring pixels have same displacement.

$$I_x(p_1)u + I_y(p_1)v = -I_t(p_1)$$

$$I_x(p_5)u + I_y(p_{25})v = -I_t(p_{25})$$

$$\Rightarrow A^T A \hat{x} = A^T b$$

$$\Rightarrow \hat{x} = (A^T A)^{-1} A^T b$$

↳ Solvable when:

$A^T A \rightarrow$  invertible

$\Leftrightarrow A^T A \rightarrow$  not too small  $\Rightarrow \lambda_1, \lambda_2 \rightarrow$  not too small.

Harris Corner Detector

↳ corner when  $\lambda_1, \lambda_2$  are big  $\Rightarrow$  Also when L-K optical flow works best.

→ Barber's pole illusion.

→ Aperture Problem

↳ Small visible image patch

↳ In which dirn is line moving?

↳ Avoid by

having patches with  
different gradients.

Smoothness: Horn-Schunck Optical Flow.

Enforce Brightness constancy,  $I_x u + I_y v + I_t = 0$

→ For every pixel,

$$\min_{u, v} [I_x u_{ij} + I_y v_{ij} + I_t]^2$$

Enforce smooth flow field,

$$\min_u (u_{ij} - u_{i+1,j})^2 \rightarrow \text{large for noisy patch}$$

Combining,

$$\min_{u, v} \sum_{i, j} \left\{ E_s(i, j) + \lambda E_d(i, j) \right\} \xrightarrow{\text{solve}} \begin{array}{l} \text{compute partial derivatives,} \\ \text{brightness constancy} \end{array} \quad \begin{array}{l} \text{Gradient descent.} \\ \text{update eqns.} \end{array}$$

$$E_s = \frac{1}{4} \left[ (u_{ij} - u_{i+1,j})^2 + (u_{ij} - u_{i,j+1})^2 \right. \\ \left. + (v_{ij} - v_{i+1,j})^2 + (v_{ij} - v_{i,j+1})^2 \right]$$

$$E_d = (I_x u_{ij} + I_y v_{ij} + I_t)^2$$

$$\frac{\partial E}{\partial u_{kl}} = 2(u_{kl} - \bar{u}_{kl}) + 2\lambda(I_x u_{kl} + I_y v_{kl} + I_t) I_x, \quad \left. \begin{array}{l} \text{if terms depend} \\ \text{on } k, l \\ \text{from} \\ \text{smoothness} \end{array} \right\}$$

$$\frac{\partial E}{\partial v_{kl}} = 2(v_{kl} - \bar{v}_{kl}) + 2\lambda(I_x u_{kl} + I_y v_{kl} + I_t) I_y \quad \left. \begin{array}{l} \text{from} \\ \text{brightness} \end{array} \right\}$$

Find extrema of  $E$ .  $\left[ \bar{u}_{ij} = \frac{1}{4} \{ u_{i+1,j}, u_{i-1,j}, u_{i,j+1}, u_{i,j-1} \} \right]$

$$\Rightarrow (1 + \lambda I_x^2) u_{kl} + \lambda I_x I_y v_{kl} = \bar{u}_{kl} - \lambda I_x I_t,$$

$$\lambda I_x I_y u_{kl} + (1 + \lambda I_y^2) v_{kl} = \bar{v}_{kl} - \lambda I_y I_t$$

$$\text{or, } A_k b$$

$$\Rightarrow x = A^{-1} b = \frac{\text{adj } A}{\det A} b$$

$$\Rightarrow \{1 + \lambda(I_x^2 + I_y^2)\} u_{kl} = (1 + \lambda I_x^2) \bar{u}_{kl} - \lambda I_x I_y \bar{v}_{kl} - \lambda I_x I_t,$$

$$\{1 + \lambda(I_x^2 + I_y^2)\} v_{kl} = (1 + \lambda I_y^2) \bar{v}_{kl} - \lambda I_x I_y \bar{u}_{kl} - \lambda I_y I_t.$$

$$\Rightarrow \boxed{\begin{aligned} \hat{u}_{kl} &= \bar{u}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_x, \\ \hat{v}_{kl} &= \bar{v}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_y \end{aligned}}$$

when  $\lambda$  is small  $\Rightarrow \lambda^{-1}$  is big

$$\Rightarrow \hat{u}_{kl} = \bar{u}_{kl}$$

$\hat{v}_{kl} = \bar{v}_{kl}$  {only smoothness}

### Horn-Schunck Optical Flow Algorithm:

① Precompute image gradients,  $I_y, I_x$ .

② Precompute temporal gradients,  $I_t$ .

③ Initialize flow field,  $u=0, v=0$

④ while not converged, compute flow field update for each pixel using  $\hat{u}_{kl}, \hat{v}_{kl}$  formula above

### Pyramidal Optimal Flow Estimation

Pyramids: 2 types: ① Low pass or Gaussian  $\hookrightarrow$  convolve image with a small gaussian kernel.

② Band pass or Laplacian  $\hookrightarrow$  subsample to get lower resolution image

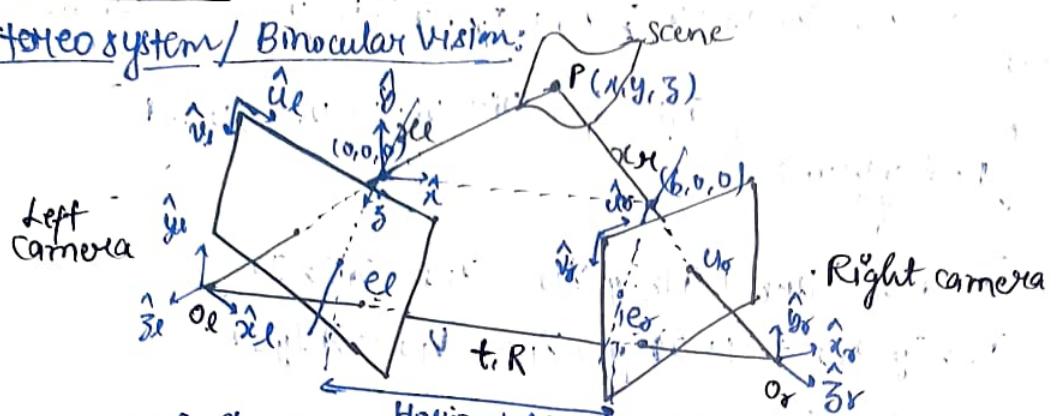
$\downarrow$  highlight edges at different spatial scales.

$\hookrightarrow$  less sensitive to illumination changes, but more noisy than using Gaussians.

$\hookrightarrow$  repeat for more levels  
A sequence of locally filtered images.

# Stereo Vision

## Stereo system / Binocular Vision:



$$u_L = fx \frac{x}{z} + 0x$$

$$v_L = fy \frac{y}{z} + 0y$$

$$u_R = fx \frac{x-b}{z} + 0x$$

$$v_R = fy \frac{y}{z} + 0y$$

Solving for  $(x, y, z)$ , if  $f, t, b, 0x, 0y$  are known.

$$x = \frac{b(u_L - 0x)}{u_L - u_R}, \quad y = \frac{bfx(v_L - 0y)}{fy(u_L - u_R)}$$

$$z = \frac{bfx}{u_L - u_R}$$

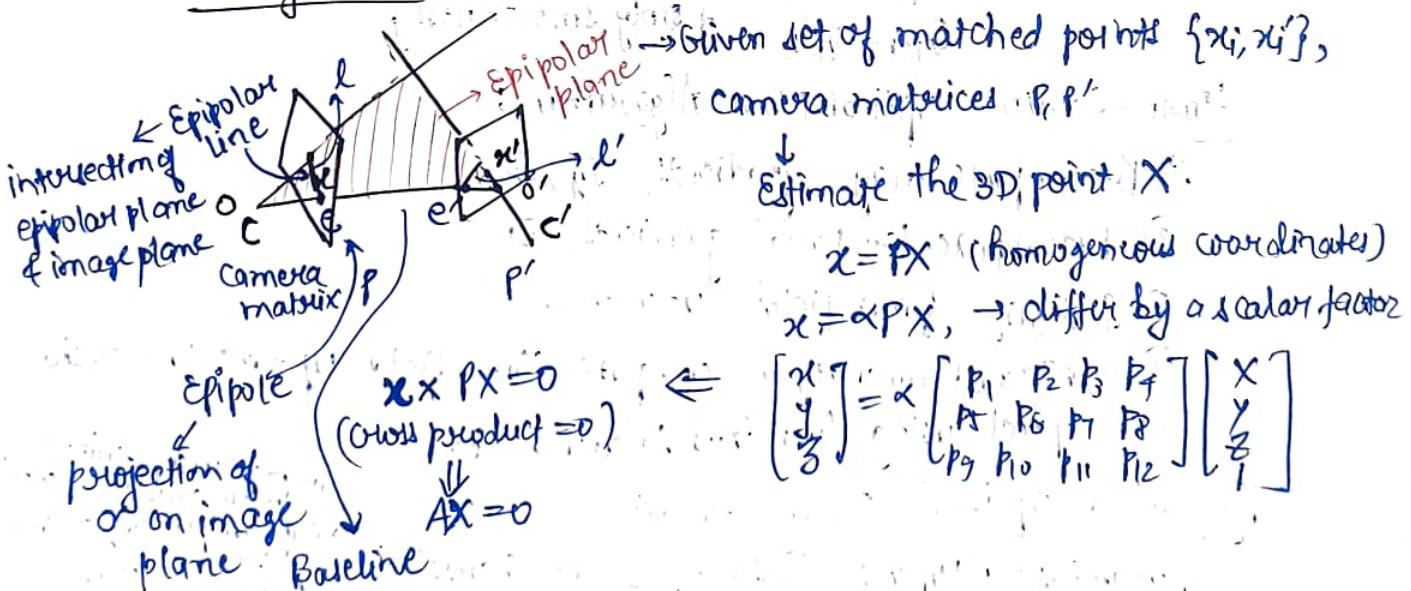
Disparity  $\propto$  Baseline (b)

Depth  $z \propto \frac{1}{\text{disparity}}$

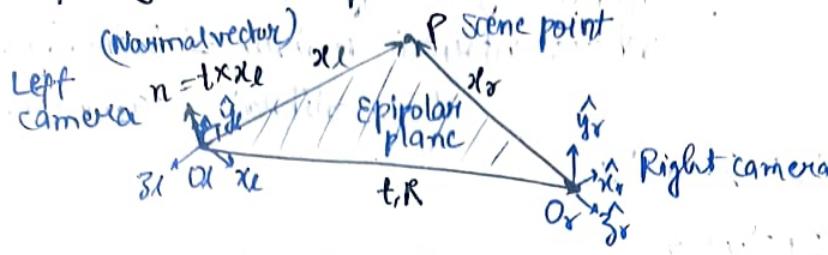
Disparity = 0

points infinitely far away from the cameras

## Triangulation:



## Computing extrinsic parameters of two cameras



$$x_L \cdot (t \times x_L) = 0$$

$$\Rightarrow [x_L \ y_L \ z_L] \begin{bmatrix} ty \ 3L - tzy_L \\ tz_L \ x_L - tx_L \\ tx_L - ty \ 3L \end{bmatrix} = 0$$

$$\Rightarrow [x_L \ y_L \ z_L] \begin{bmatrix} 0 & -tz & ty \\ tz & 0 & -tx \\ -ty & tx & 0 \end{bmatrix} \begin{bmatrix} x_L \\ y_L \\ z_L \end{bmatrix} = 0$$

Substitute  $x_L = Rx + t$ .

$t^3x_L$ : position of R.C. in LC's frame

$R_{3x3}$ : orientation of L.C. in R.C.'s frame

$$\Rightarrow [x_L \ y_L \ z_L] \left( \begin{bmatrix} 0 & -tz & ty \\ tz & 0 & -tx \\ -ty & tx & 0 \end{bmatrix} \begin{bmatrix} x_L & y_L & z_L \\ x_L & y_L & z_L \\ x_L & y_L & z_L \end{bmatrix} \begin{bmatrix} x_L \\ y_L \\ z_L \end{bmatrix} + \begin{bmatrix} 0 & -tz & ty \\ tz & 0 & -tx \\ -ty & tx & 0 \end{bmatrix} \begin{bmatrix} tx \\ ty \\ tz \end{bmatrix} \right) = 0$$

$$E = \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} \quad t \times t = 0$$

$$\Rightarrow [x_L \ y_L \ z_L] E \begin{bmatrix} x_L \\ y_L \\ z_L \end{bmatrix} = 0 \quad \text{, where } E = T \times R \rightarrow \text{essential matrix}$$

$$\Rightarrow x_L^T E x_L = 0 \quad \rightarrow \text{Relative 3D position of scene w.r.t. LC to its 3D position w.r.t. RC.}$$

→ Potential matches to  $x$  lie on the epipolar line  $l'$ .

→ Epipolar constraint reduces search to a single line.

$$Ex = l'$$

$$x'^T Ex = 0$$

$\because x^T Ex = 0, Ex = l'$  → Maps a point to a line

$$x' = R(x-t) \quad (\text{Rigid motion}) \quad \rightarrow \text{different from homography.}$$

$$x, t, x' \text{ are coplanar} \Rightarrow x^T(t \times x) = 0 \quad \Rightarrow (x^T R)(t \times x) = 0$$

### Properties of Essential Matrices

Longuet-Higgins equation:  $x'^T Ex = 0$

Epipolar lines:  $x^T l = 0, x'^T l' = 0$

$$l' = Ex, l = E^T x'$$

Epipoles:  $e^T E = 0, Ee = 0$

$$\begin{cases} x^T Ex = 0 \\ E = (R[t \times]) \end{cases}$$

Assumption: 2D points expressed in camera cs.

(i.e., intrinsic matrices are identities.)

Fundamental Matrix: Generalization of essential matrix, where intrinsic matrix  $\neq$  identity.

$$\hat{x}'^T E \hat{x} = 0$$

↳ uncalibrated camera.

$$\hat{x}' = K'^{-1}x', \hat{x} = K^{-1}x$$

$$\Rightarrow K'^{-T} E K^{-1} x = 0$$

$$\Rightarrow x'^T (K'^{-T} E K^{-1}) x = 0$$

$$\Rightarrow \boxed{x'^T F x = 0}$$

→ Properties of essential matrix applicable.

$$F = K'^{-T} E K^{-1}$$

$$= K'^{-T} [f_x] R K^{-1}$$

↳ independent both intrinsic & extrinsic parameters.

Solve for F:  $x_m^T F x_m = 0$ , → for M matched image points.

$$\Rightarrow \begin{bmatrix} x_m' & y_m' & 1 \end{bmatrix} = \begin{bmatrix} f_1 & f_2 & f_3 \\ f_4 & f_5 & f_6 \\ f_7 & f_8 & f_9 \end{bmatrix} \begin{bmatrix} x_m \\ y_m \\ 1 \end{bmatrix} = 0.$$

$[F \rightarrow 3 \times 3]$

$$\Rightarrow x_m x_m f_1 + x_m y_m f_2 + x_m f_3 + y_m x_m f_4 + y_m y_m f_5 + y_m f_6 + x_m' f_7 + y_m' f_8 + f_9 = 0 \Rightarrow \text{One correspondence gives one equation.}$$

setup a homogeneous linear system with 9 unknowns.

Each point pair contributes only one scalar equation.

Need at least 8 points. → 8-Point Algorithm.

$$\Rightarrow Ax = 0 \rightarrow \text{solve using SVD}$$

8-point algorithm:

① (Normalize points)

② construct the  $M \times 9$  matrix A.

③ Find the SVD of A.

④ Entries of F are the elements of column of V corresponding to the least singular values.

⑤ (De-normalize F)

Enforcing Rank constraints: Given matrix F, find  $F'$  of rank K that is closest to F.

$$\min_{F'} \|F - F'\|^2 \Rightarrow \text{compute SVD of } F \Rightarrow F = U \Sigma V^T$$

→ Form  $\Sigma'$  by replacing all but largest singular value in  $\Sigma$  with 0.

$$\Rightarrow \therefore F' = U \Sigma' V^T.$$

## Image segmentation

- Partition an image into a collection of set of pixels (segments).
- Meaningful regions, shapes.
- Techniques: Thresholding; Region-based methods (region growing), clustering (k-means), graph-based methods (graph-cut, random), shape-based methods (level set, active contours), energy minimization methods (MRF), Machine Learning based methods.

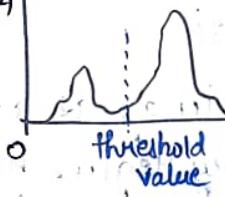
Image Binarization: → One global threshold,  $T$  for mapping scalar image  $I$  to binary image:

$$J(x,y) = \begin{cases} 0, & \text{if } I(x,y) < T \\ 1, & \text{otherwise.} \end{cases}$$

- Global threshold can be identified by an optimization strategy aiming at creating connected regions and at reducing the no. of small-sized regions, called artifacts.

### Thresholding:

#### Histogram:



#### Difficulties:

- Valley may be broad.
- No. of minima.
- Noise; no visible valley.
- Multi-modal histogram.

Otsu Thresholding: Uses gray-value histogram to provide best threshold such that overlap b/w 2 classes (object and background pixels) is minimized.

$$\sigma_b^2 = P_1(\mu_1 - \bar{\mu})^2 + P_2(\mu_2 - \bar{\mu})^2 = \bar{\mu}P_2(\mu_1 - \mu_2)^2$$

$\sigma_b^2$  : b/w-class variance

$P_1, P_2$  : class probabilities

$\mu_i$  : means of object and bg. classes.

Algorithm: Compute histogram  $H_I$  for  $i=0, \dots, G_{\max}$ .

Let  $T_0$ : increment for potential thresh.,  $u=T_0$ ,  $T=u$ ,  $s_{\max}=0$ .  
while  $u < G_{\max}$  do

compute  $C_I(u)$ ,  $\mu_i(u)$  for  $i=1, 2$ .

compute  $\sigma_b^2 = [C_I(u) - C_I(u)] [\mu_1(u) - \mu_2(u)]^2$

if  $\sigma_b^2(u) > s_{\max}$  then

$s_{\max} = \sigma_b^2$ ,  $T=u$

set  $u=u+T_0$ .

$$\# \quad P = \sum_{i=0}^u p(i), \quad P_1 = \sum_{i=u+1}^{G_{\max}} p(i), \quad M_1 = \sum_{i=0}^u i p(i)/P, \quad M_2 = \sum_{i=u+1}^{G_{\max}} i p(i)/P_1$$

$C_I$ : Relative cumulative histogram of image  $I$ .

$$P \approx C_I(u), \quad P_2 \approx 1 - C_I(u)$$

$u$ : chosen threshold.

- ↳ Performs well when histogram has a bimodal distribution with a deep and sharp valley b/w 2 peaks.
- ↳ Performs badly in case of heavy noise, small objects size, inhomogeneous lighting, larger inter-class & intra-class variance.

### Region-based Segmentation

Region Growing: Full segment generation in grayscale or color images, start at one seed pixel  $(x, y, I(x, y))$ , and add recursively adjacent pixels, that satisfy a "similarity criterion" with pixels contained in the so far grown region around seed pixel.

### Region Splitting and Merging Segmentation

Region Splitting: Starts with the whole image as a single region and subdivides it into subsidiary regions recursively while a condition of homogeneity is not satisfied.

Region merging: Opposite to splitting, to avoid over-segmentation.

- ↳ Start with a small region, and merge the regions that have similar characteristics (gray level, variance, etc.).

### Clustering based Segmentation

- ↳ Organizing data into classes such that - high intra-class similarity, low inter-class similarity.
- ↳ Finding the class label of the no. of classes directly from the data.
- ↳ Cluster by features.
- ↳ Mean shift segmentation
- ↳ Partitioning image into meaningful regions or objects by grouping pixels with similar features, like color.