



Computer Vision

1. Introduction

→ Computer programs that can interpret images.
Understanding the representation of image.

→ Applications:

- Images & movies

ex- Surveillance

Building 3D representations → Earth, Movies

Motion capture → Pirates of Caribbean, CGI

- OCR and Face Recognition.

↓
Scanner → Adobe, Supermarket

- Object Recognition

- Special Effects & 3D Modeling

- Smart Cars

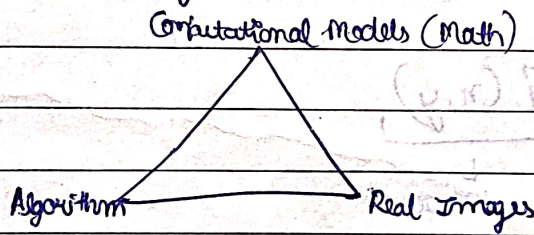
- Sports → Cricket, Soccer.

- Vision Based Interaction → MS Kinect, Nintendo Wii

- Security and surveillance.

Notes

Image Processing & computer vision is different.



2.

Image Processing for CV

Image is a function, f on I , from R^2 to R

Blurry Image = Smooth function

$f(x, y)$ gives the intensity or value at position (x, y)

Practically define image over a rectangle $f: [a, b] \times [c, d] \rightarrow [\min, \max]$

with a finite range.

Notes 0 to 255 is usual measurement (0 to 1 should be better)

for color image

$$f(x,y) = \begin{bmatrix} r(x,y) \\ g(x,y) \\ b(x,y) \end{bmatrix}$$

• Image Blending

$$g(x) = (1-\alpha) f_0(x) + \alpha f_1(x) + \gamma$$

↑
usually zero

* pixel → picture element.

$$dst = cv2.addWeighted(img1, 0.7, img2, 0.3, 0)$$

→ Digital Images:

In computer vision we typically operate on digital images.

Sample → the 2D space on a regular grid.

Quantize each sample (round to "nearest integer")

Images, they are represented as matrix of integers.

→ Noise in Images

$$I'(x,y) = I(x,y) + \underbrace{n(x,y)}_{\text{noise}}$$

Types of Noise

- Salt & pepper: random occurrences of black & white pixels
- Impulse: random occurrences of white pixels
- Gaussian: variations in intensity drawn from a Gaussian normal distribution

↓

* noise = randn(size(img)) * sigma ← better depends on range of image

output = im + noise

Image Filtering

used to

→ Enhance image

→ Extract information

→ Detect patterns

→ Removing noise from an image:
(Moving Average) 1-D

Approach 1: Replace each pixel with an average of all the values in its neighbourhood.

• Assumptions

1. The "true" value of pixels are similar to the true value of pixels nearby.
2. The noise added to each pixel is done independently.

Approach 2: (Weighted Moving image) 1-D
Can add weights to our moving average

So, we use an odd number of ^{non-uniform} weighted mask

eg. [1, 4, 6, 4, 1] to all nearby neighbourhood

For 2-D

Approach 1: Moving Average

Just take some ~~2x2~~ or 3x3 matrix and average it.

* Co-relation filtering - uniform weights

Say the averaging window size is $(2k+1) \times (2k+1)$

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

Nested loop

↓
& averaging it.

* Correlating filtering - non-uniform weights

$$G[i, j] = \sum_{u=-K}^K \sum_{v=-K}^K H[u, v] F[i+u, j+v]$$

This is called cross-correlation, denoted $G = H \otimes F$

This $H[u, v]$ is called filter "kernel" or "mask" which is just the matrix of weights in the linear combination.

So, what makes a good kernel?

Approach 2: Gaussian Filter

$$H(u, v) \rightarrow \begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 2 & 4 & 2 \\ \hline 1 & 2 & 1 \\ \hline \end{array}$$

$$h(u, v) = \frac{1}{2\pi\sigma^2} e^{-\frac{u^2+v^2}{2\sigma^2}}$$

This kernel is approximation of Gaussian function

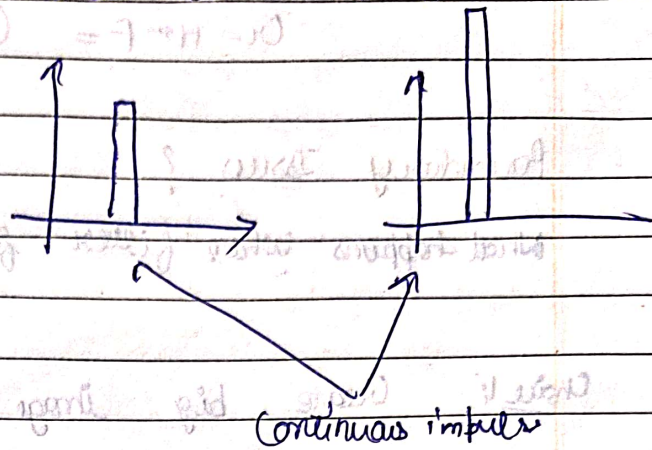
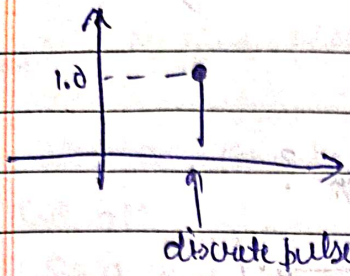
Note: Variance (σ^2) or standard deviation (σ) determines extent of smoothing

Also, Size of Kernel (u, v) is not variance

⇒ Impulse Noise

In discrete world, impulse is just a value of 1 at a single location.

while in continuous world, impulse is "idealized funcⁿ" that is very narrow and very tall so that it has a unit area.



• Convolution

$$G_i(i, j) = \sum_{u=-K}^K \sum_{v=-K}^K H[u, v] F[i-u, j-v]$$

$$G_i = H * F$$

Also, $\text{Img} * \text{Impulse} = \text{Img}$

→ Properties of Convolution:

a) Linear & Shift invariant

b) $f * g = g * f$ (Commutative)

c) $(f * g) * h = f * (g * h)$ (Associative)

d) $f * e = f$ (Identity)

e) Differentiation $\frac{\partial}{\partial x} (f * g) = \frac{\partial f}{\partial x} * g$

A Big Note: What would be the computational complexity?

$W * W * N * N$ to get the convolution
 Filter Image

so, if some filter is separable we can create convolution by $W * N * N + W * N * N = 2 * W * N * N \ll W^2 * N^2$

It can be done because

$$C_{in} = H \times F = (C \times R) \times F = C \times (R \times F)$$

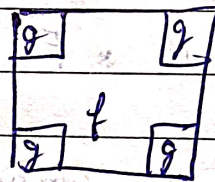
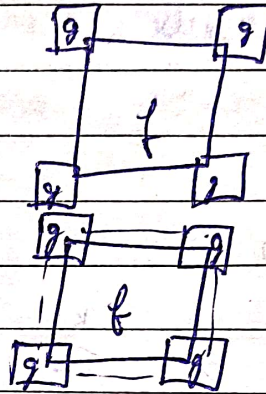
Boundary Issues ?

What happens when filter falls off the edge ?

Choice 1: Create big image

Choice 2: Create same image

Choice 3: Create smaller image



Methods ? (In opencv) [For same image]

i) cv2.BORDER_CONSTANT \rightarrow just clips out extra & replace with some color

ii) cv2.BORDER_WRAP \rightarrow Circular wrapping

iii) cv2.BORDER_REPLICATE \rightarrow last element is replicated throughout

iv) cv2.BORDER_REFLECT \rightarrow Border will be mirror reflection of the border elements.

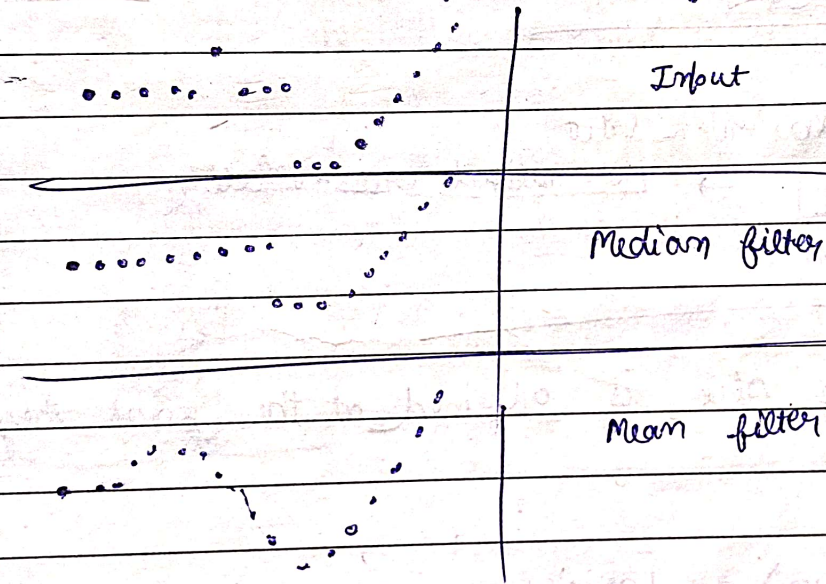
No obvious edges

Now, talking about filters for noises like salt & pepper noise, impulse noise

Approach 3: Median Filter

- Select a kernel (odd \times odd = odd)
- Sort all the numbers
- Select mid element
- Replace it

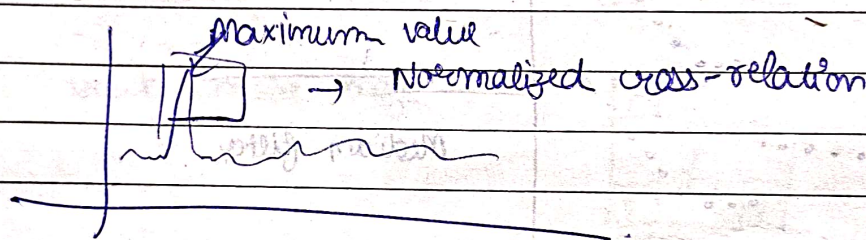
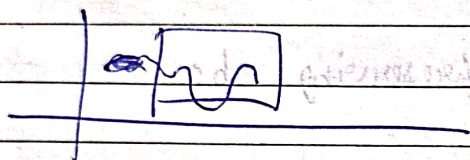
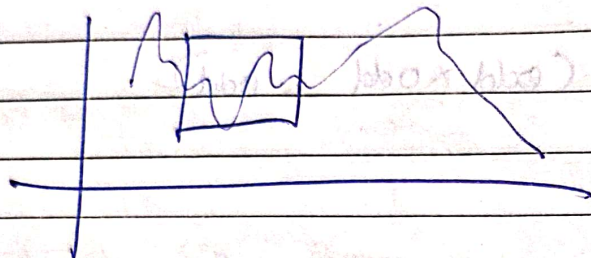
It is also called edge preserving bc



Now, we move on to finding properties of pixels locally in the image

→ Template Matching :

- 1D (nx) correlation



Why maximum value is attained at the point they match completely?

$$(+ve) \times (+ve) = \text{positive}$$

$$(-ve) \times (-ve) = (+ve)$$

Well, the filter max out each signal when they match exactly.

★ We can now use this trick to find objects in the image (Template Matching)

⇒ Here, this doesn't work well for different scales & orientations.

So, it's just an applⁿ of using filters as template.

⇒ Edge Detection: Gradients

Edges in an image is important, they can easily convey information with reduced images.

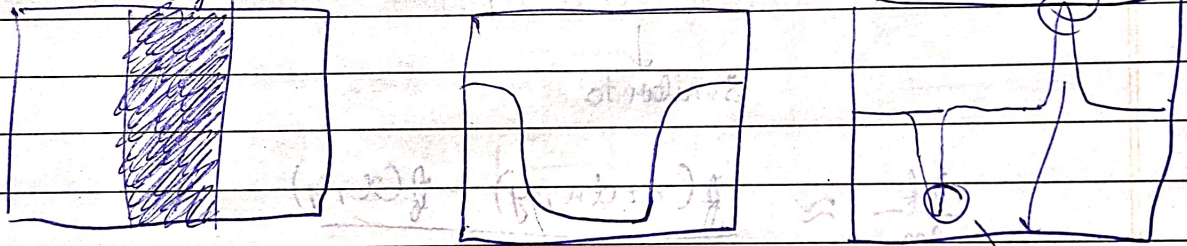
Matrix $\rightarrow f(x, y)$ $\xrightarrow{\text{reduce}}$ reduced set of pixels or curves that are important elements of the picture.

Now, edges look like steep cliffs / change, so look for a neighborhood with strong signs of change.

But, there are some ambiguity

- neighborhood size
- how to detect change

Now, change is derivative (as an edge is a place of rapid change in the image intensity function)



Edges correspond to extrema of derivative

→ To find these peaks (extrema of derivative), we run it through some operators (differential operators to be exact)

- Differential operators - when applied to some image returns some derivatives (duh?)

- Model these "operators" as masks / kernels that compute the image gradient function.

- Threshold the this gradient function to select the edge pixels.

But, what is gradient?

Gradient is a vector that's made up of the partial derivatives

$$\text{The gradient of an image: } \nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$$

In other words, the gradient points in the direction of most rapid increase in intensity, & magnitude of vector is how much it's changing as a function of a unit step in that direction.

For computers, we calculate discrete gradient

$$\frac{\partial f(x, y)}{\partial x} \approx \frac{f(x+1, y) - f(x, y)}{1}$$

↓
Similar to

$$\frac{\partial f}{\partial x} \approx \frac{f(x+dx, y) - f(x, y)}{dx}$$

Some operators

1. Sobel operator

$$S_x = \frac{1}{8} * \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

$$S_y = \frac{1}{8} * \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

It is a joint gaussian smoothing plus differentiation operator, & it is more resistant to noise.

$$\nabla I = \begin{bmatrix} g_x & g_y \end{bmatrix}^T$$

→ Gradients to edges:

Primary edge detection steps:

1. Smoothing derivatives to suppress noise and compute gradient.
2. Threshold to find "regions" of "significant" gradient
3. "Thin" to get localized edge pixels
4. and link or connect edge pixels.

So, to do that we use canny edge detection algorithm,

It's stages are:

1. Filter image with derivative of Gaussian
2. Find magnitude & orientation of gradient

$$\text{Edge-gradient } (G) = \sqrt{G_x^2 + G_y^2}$$

$$\text{Angle } (\theta) = \tan^{-1}\left(\frac{G_y}{G_x}\right)$$

3. Non-maximum suppression

Thin multi-pixel wide "ridges" down to single pixel width.
In short, the result you get is binary image with "thin edges".

4. Linking and Thresholding (Hysteresis):

- Define two thresholds: low and high
- Use the high threshold to start edge curves and the low threshold to continue them

* Some concerns

1. What's the value of σ in gaussian filter

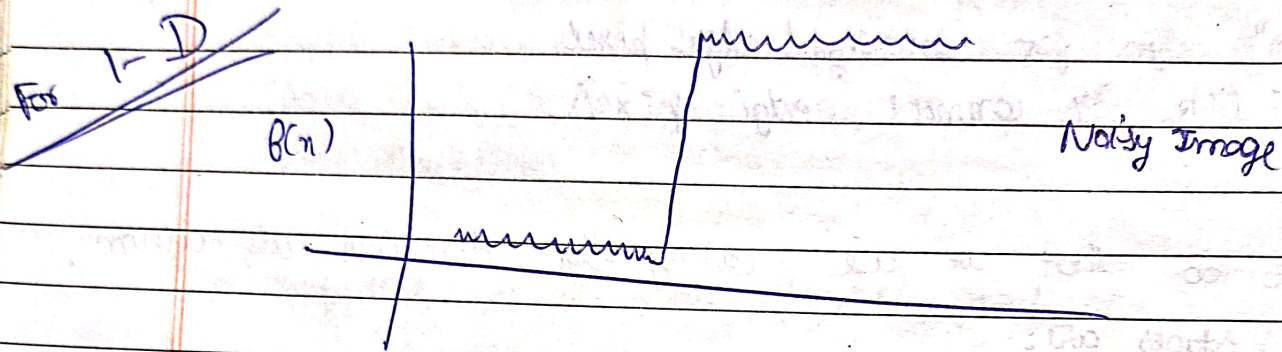
Canny with $\sigma=1$ (fine features) (detailed edge image)
Canny with $\sigma=2$ (less features) (large scale edges)

2. What's the kernel size to be taken?

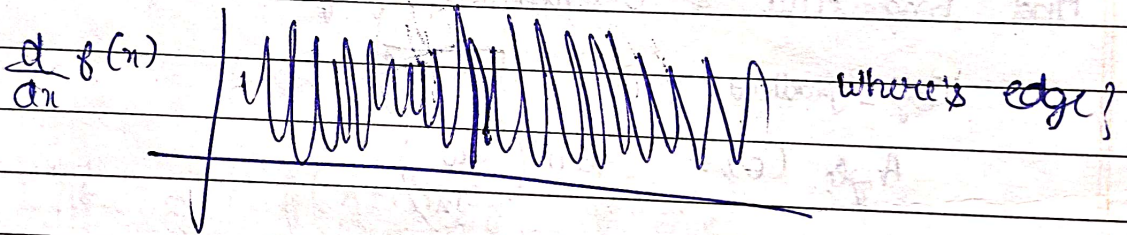
Mathematics

$f \rightarrow$ Image funⁿ

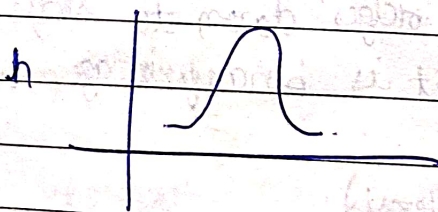
$h \rightarrow$ filter (Gaussian)



Apply derivative operator

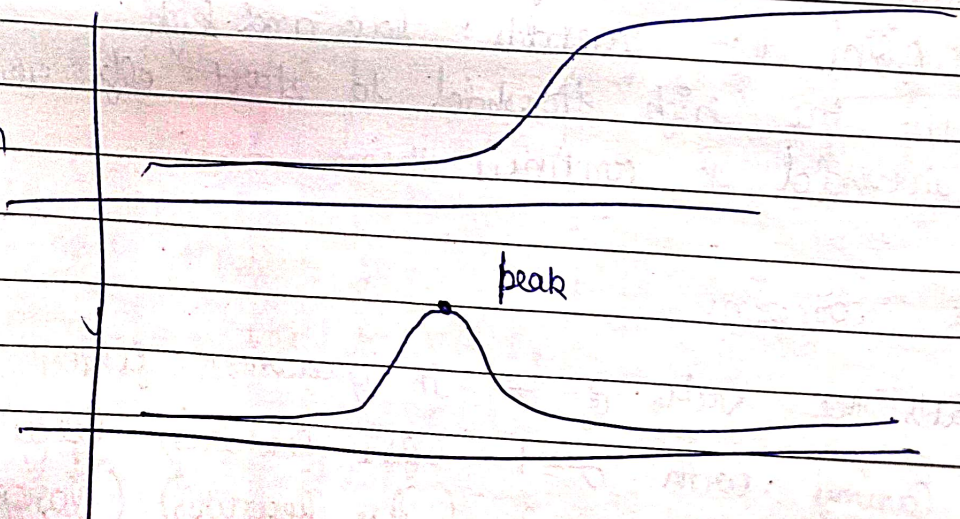


So, we need to remove noise

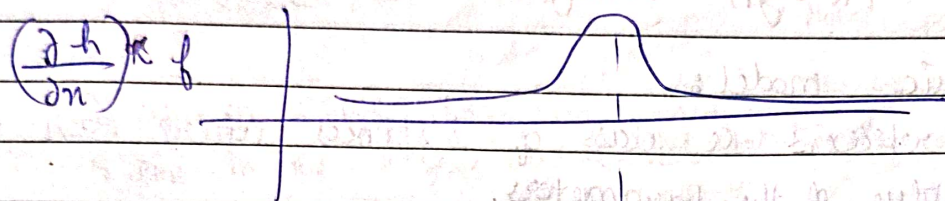
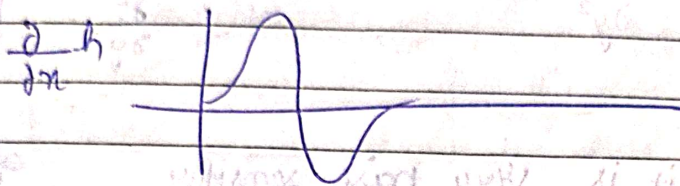


hac f
(conduction)

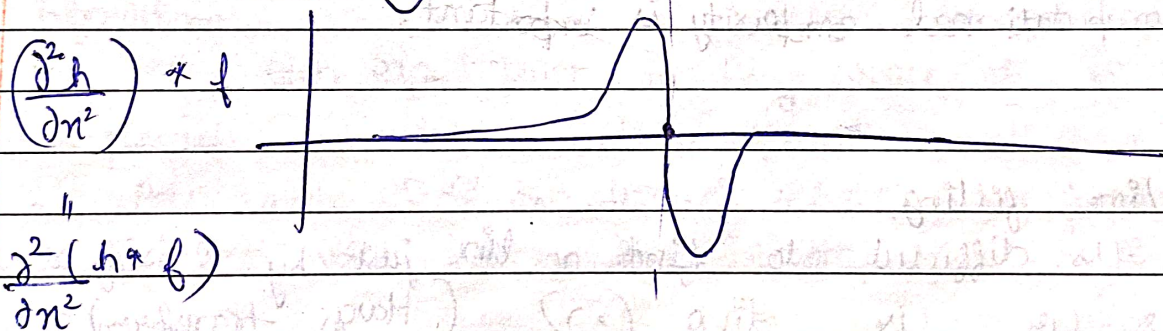
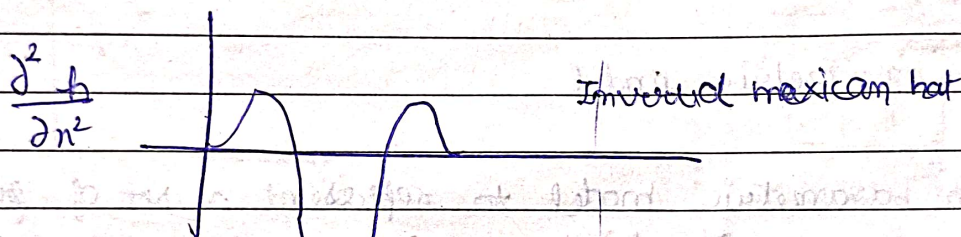
$\frac{d(hac f)}{dn}$



Also, $\frac{\partial (h * f)}{\partial n} = \left(\frac{\partial h}{\partial n}\right) * f$



Now, to find peaks, we take more derivatives



(It's a single 2D edge detection filter)

For 2-D It's a little bit harder, as there's more than one direction to take our derivative.

Gaussian $h_\sigma(u, v) = \frac{1}{2\pi\sigma^2} e^{-\frac{u^2+v^2}{2\sigma^2}}$

$\frac{\partial h}{\partial n}$ becomes tedious

$\frac{\partial^2 h}{\partial (n, y)^2}$ is Mexican Hat.

$$\nabla^2 h = I_{xx} + I_{yy}$$

$$= \left([1 \ -2 \ 1] + \begin{bmatrix} -1 \\ -2 \\ 1 \end{bmatrix} \right) * I$$

$$= \begin{bmatrix} 0 & -4 & 0 \\ 0 & -4 & 0 \end{bmatrix} * I$$

Laplacian ←

So, we use something called Laplacian operator

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

Laplacian operator

where $\frac{\partial^2}{\partial x^2} = [1 \ -2 \ 1]$
 $\frac{\partial^2}{\partial y^2} = \begin{bmatrix} -1 \\ -2 \\ 1 \end{bmatrix}$

Central diff approx to second derivative

Something required as it is very noise sensitive.
 Fourier Transform: Lines

⇒ Parametric models:

It can represent a class of instances where each is defined by a value of the parameters.

ex- lines, circles, parameterized template

→ Fitting a parametric model

- Choose a parametric model to represent a set of features.
- Membership criterion is not local.
- Computational complexity is important.

I Line fitting

It is difficult to find a line just by looking at the points.
 so we use voting (:) (Hough transform)

It is a general technique where we let the features vote for all models that are compatible with it.

1. Cycle through features, each casting votes for model parameters.
2. Look for model parameters that receive a lot of votes.

→ Hough space

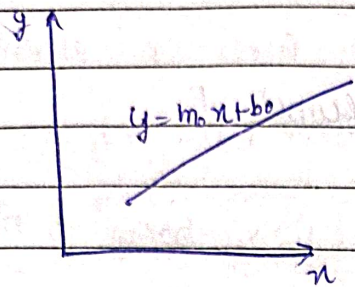
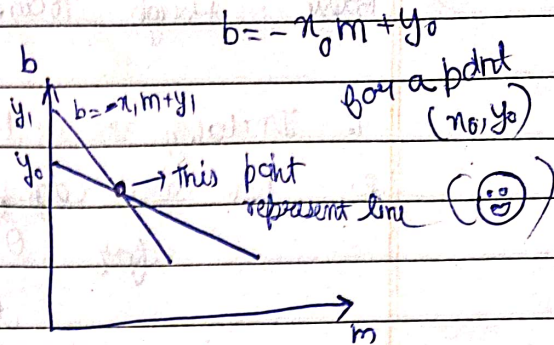
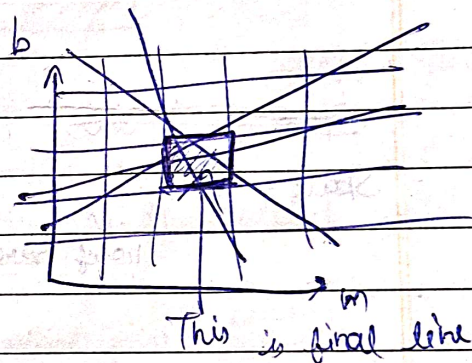
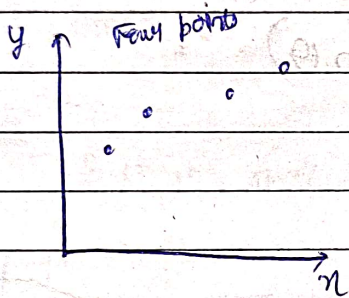


Image Space



Hough (parameter) space

A point in image space is line in Hough space
 A line in the image correspond to Hough Space point

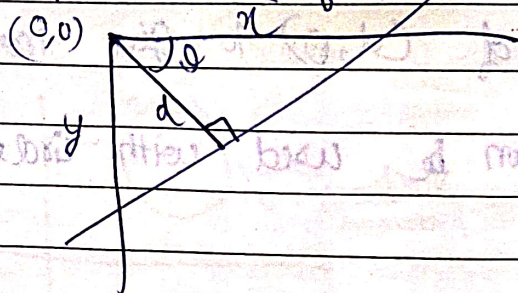


Algorithm:

1. let each edge points in image space vote for a set of possible parameters in Hough space
2. Accumulate votes in discrete set of bins; parameters with the most votes indicate line in image space.

• A little getcha

if line is vertical $m \rightarrow \infty$, so we use polar representation of a line. to prevent any mathematical error.



d : max distance from line to origin

θ : angle the line makes with x -axis

$$x \cos \theta + y \sin \theta = d$$

- The polar representation of image causes the Hough space to have sinusoidal waves for each point.

→ Basic Hough transform algorithm

1. Initialize $H[d, \theta] = 0$

2. For each edge point in image funcⁿ F
for $\theta = 0$ to π (180)

$$d = x \cos \theta + y \sin \theta$$

$$H[d, \theta]++$$

3. Find the value(s) of (d, θ) where $H[d, \theta]$ is max.

4. The detected line is given by
 $d = x \cos \theta - y \sin \theta$

~~Space = $O(n^2)$~~

$$\text{Time} = O(n^2)$$

Space = $O(K^2)$

no. of parameters

There are some extensions to the above algo., that helps with noise, line segments etc.

I Using the gradient (knowing the alignment at (x, y) in image space

instead of using $\theta = 0$ to π

use, $\theta = \text{grad.}$ or $(-35$ to $55)$

II Give more votes for stronger edges.

III Change the sampling of (d, θ) to give more/less resolution.

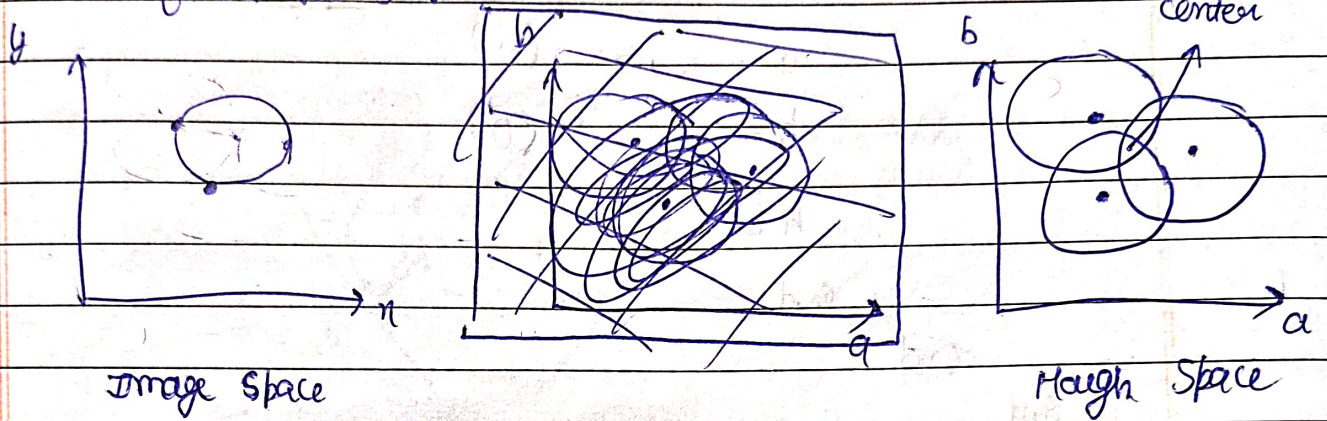
IV The same procedure can be used with circles, squares or any other shape.

II Circle

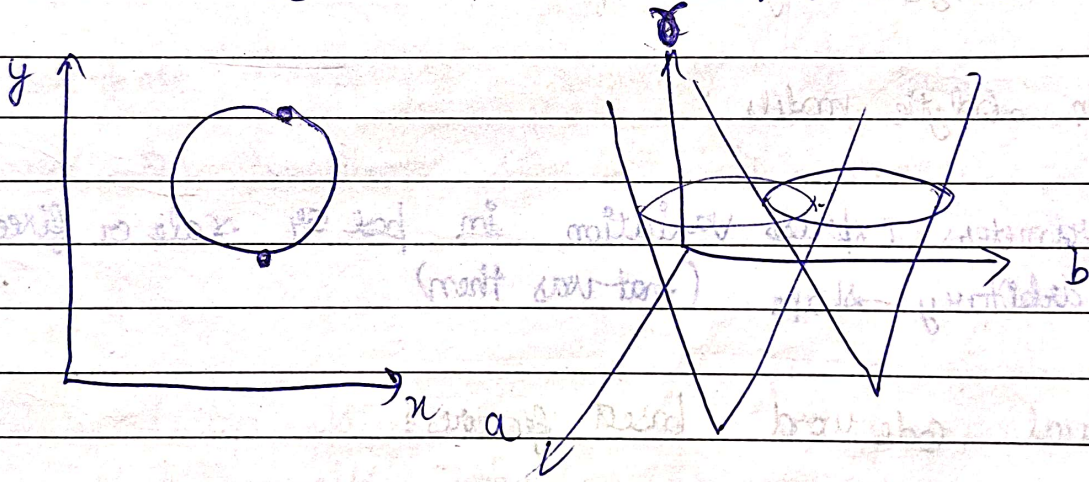
Consider a circle with center (a, b) & radius r

$$(x_i - a)^2 + (y_i - b)^2 = r^2$$

Method 1: For a fixed radius r



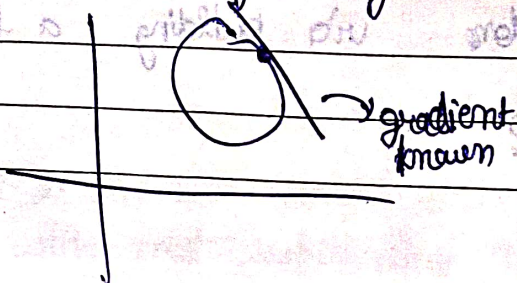
Method 2: But in real world, we don't know the radius so



The Hough space increases drastically. (RANSAC solves it)

Method 3: Unknown r , with gradients (reduces Hough space)

Vote only along a line or circle.



i) If r is known
line

ii) If r is unknown
circle

→ Algorithm

1. For every edge pixel (x, y)

For each possible radius value r :

For each possible gradient direction θ :

%% or use estimated gradient

$$a = x + r \cos(\theta)$$

$$b = y + r \sin(\theta)$$

$$H(a, b, r) += 1$$

end

end

end

III. Generalized Hough Transform

⇒ Non analytic models:

Parameters express variation in pose or scale or fixed but arbitrary shape (that was then)

⇒ Visual code-word based features:

(Not edges) but detected templates learned from models (this is "now")

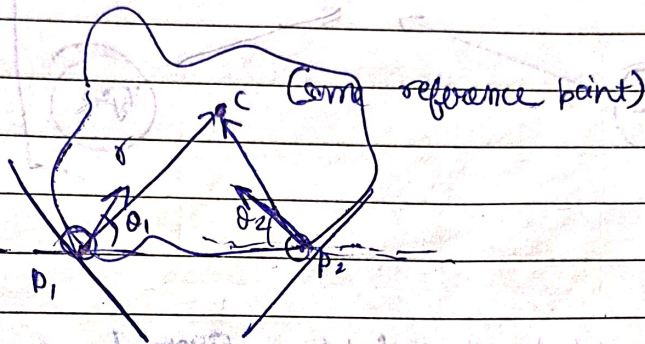
Moving back to Generalized Hough transform

In this, voting is done via building a hough table.
How to build Hough table?



Training

1. At each boundary point, compute displacement vector $\delta = c - p_i$.
2. Measure the gradient angle θ at the boundary point.
3. Store that displacement in a table indexed by θ .



Recognition

1. At each boundary point, measure the gradient angle θ .
2. Look up all displacements in θ displacement table.
3. Vote for a center at each displacement.

Algorithm

If orientation is known:
do recognition point

else

for each edge point

for each possible master θ^*

compute gradient direction θ

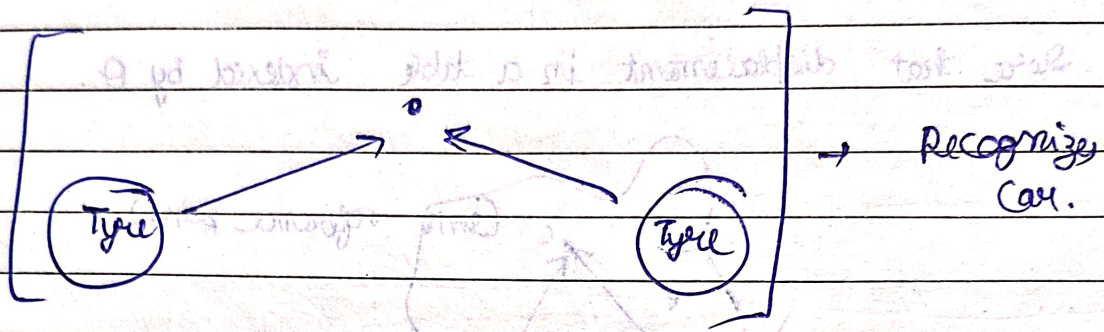
New $\theta' = \theta - \theta^*$

for θ' retrieve displacement vectors to vote for reference point.

Peak in this Hough space (now X, Y, θ^*) is reference point with most supporting edges.

• Modern Approach:

Instead of using indexed displacements by gradient orientation, index by "visual codeword"



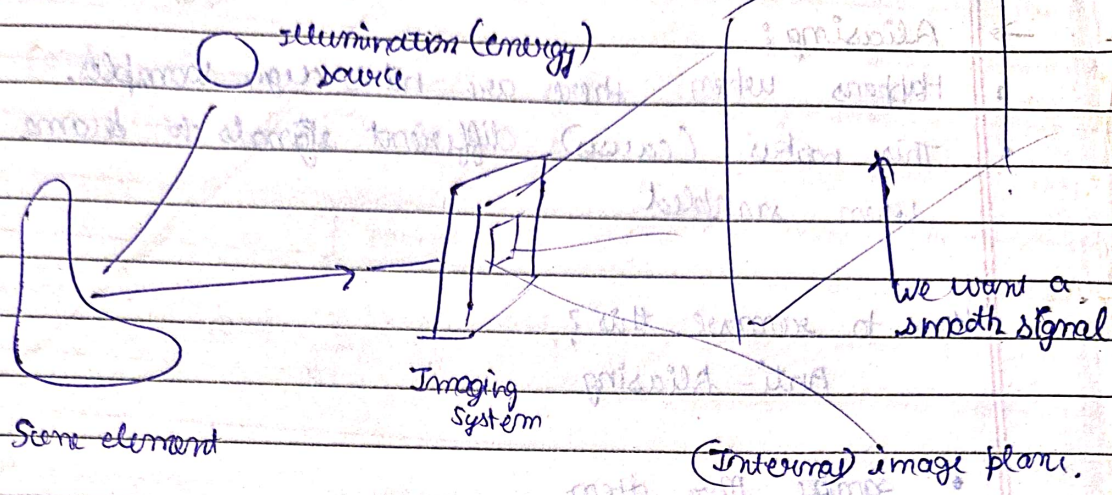
Training

1. Build codebook of patches around extracted (interest points) using clustering.
2. Map the patch around each "interest point" to closest codebook entry.
3. For each codebook entry, store all displacements relative to object center.

Recognition

1. Find all the "interest point" patches & identify visual codeword.
2. Look up all the displacements in codebook.
3. Vote for center at each displacement.

Aliasing

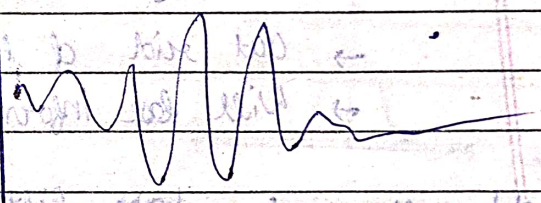


→ Sampled representations

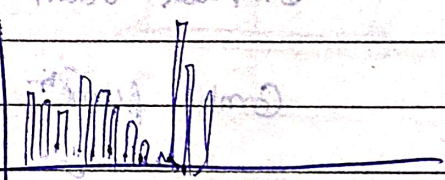
• Sampling → How to store & compute with continuous funcⁿ.

* If we do ~~undersampling~~, information is lost.

But, it was indistinguishable from lower frequency.

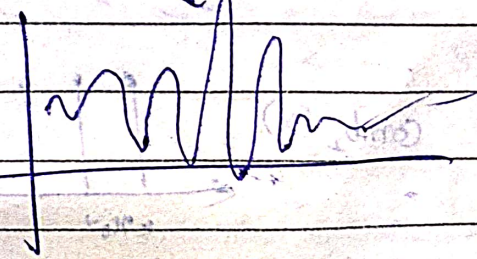


↓ Sampling

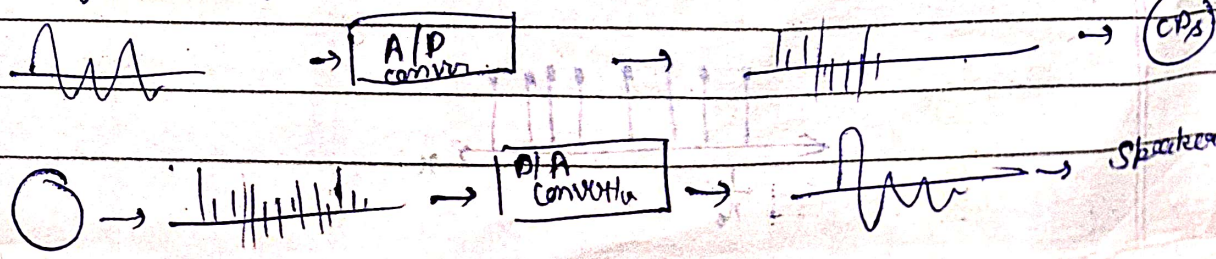


Goal → "guess" what this funcⁿ did in between.

↓ Reconstruction



• Sampling in digital audio



→ Aliasing:

- Happens when there are not enough samples.
- This makes (causes) different signals to become indistinguishable when sampled

How to remove this?
Anti-Aliasing

- Sample More often
→ But this can't go on forever
- Make the signal less "wiggly"
→ Out side of higher frequencies
→ Will lose information

method 1 use low pass filters

(image is stored already)

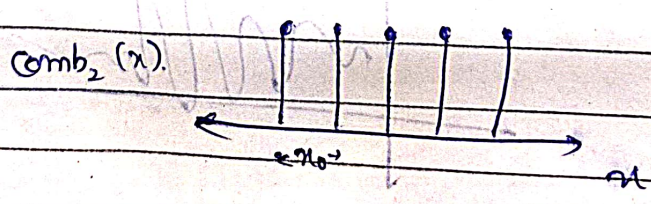
method 2 Impulse Train

Comb funcⁿ (in 1D)

$$= \sum_{k=-\infty}^{\infty} \delta[x - km]$$

$$\delta(x) = 1 \quad \text{if } x=0$$

$$\delta(x) = 0 \quad \text{otherwise}$$



Fourier Transform (Scaling property)



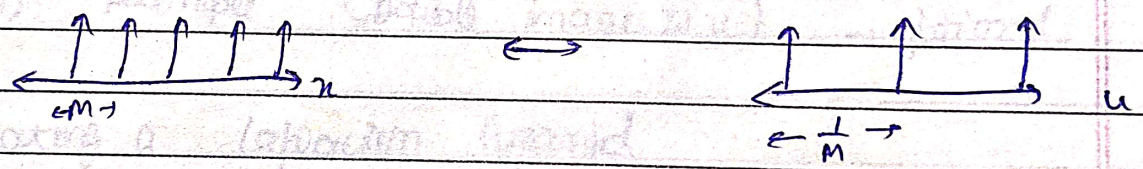
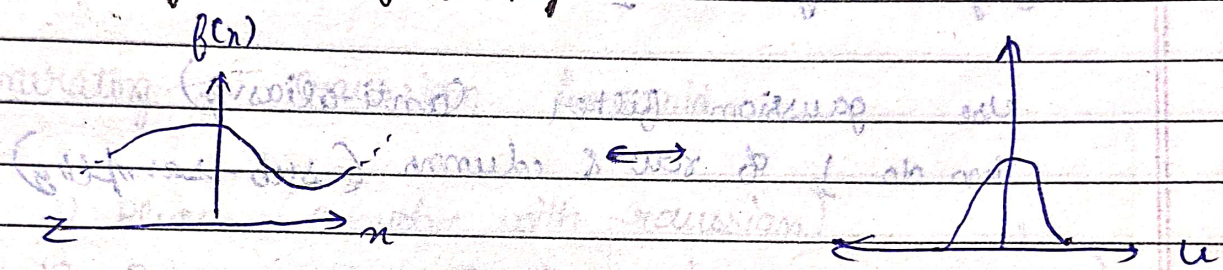
Impulses in 2D

$$\text{comb}_{M,N}(x,y) = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} \delta(x - kM, y - lN)$$

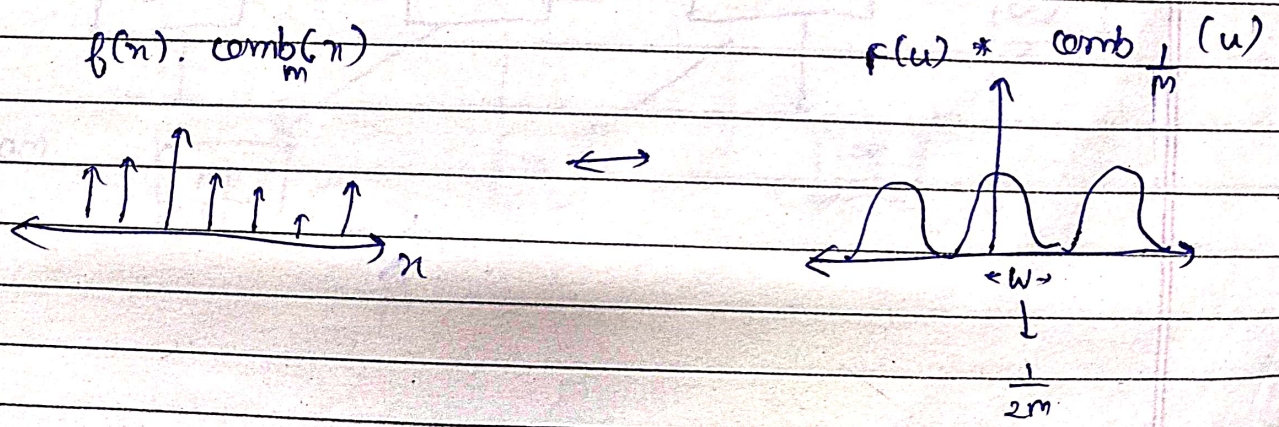
Fourier

$$\text{comb}_{\frac{1}{M}, \frac{1}{N}}(u,v)$$

To Sampling law frequency



multiply (sample)



If there is no overlap, $\Delta u < \frac{1}{2m}$, the original signal can be recovered from its sample by low pass filtering.

II. Sampling high frequency:

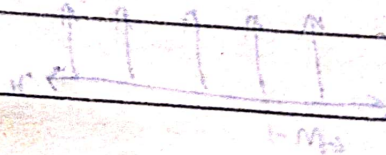
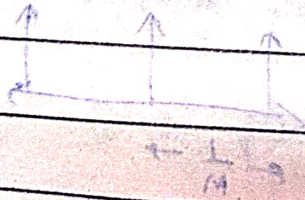
- Apply an anti-aliasing filter
- Then do comp funcⁿ. ~~Comp~~

→ Aliasing in Images:

Resizing image in $\frac{1}{2}$

Use gaussian filter (anti-aliasing)

Then do $\frac{1}{2}$ of row & column. (sub-sampling)



(w)

(w)

* Comp

(w)

(w)

(w)

⇒ Image Pyramids

• there are two kinds of image pyramids

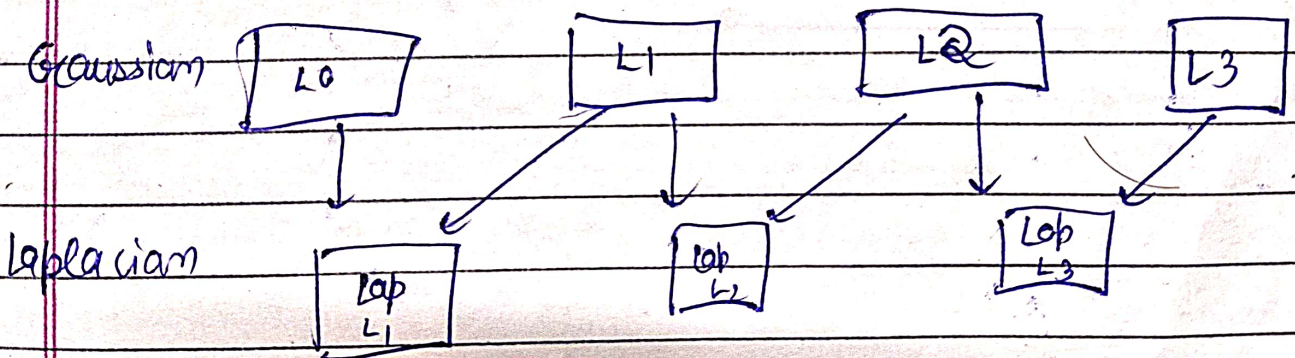
- i) Gaussian Pyramid
- ii) Laplacian Pyramids

We use these pyramids, when we need to work with images of different resolution of the same image.

• Generating a Gaussian Pyramid

- i) Blur (convolve with Gaussian)
- ii) Downsample (reduce image size by half)
- iii) Upsample (double image size) → optional

• Generating a Laplacian Pyramid



→ Reduce & Expand :

$$\left[\begin{array}{cccc} 1 & 4 & 6 & 4 & 1 \end{array} \right] / 16$$

→ Reduce } Repeatable
 5-tap filter

$$\left[\begin{array}{ccccc} \frac{1}{8} & \frac{1}{2} & \frac{3}{4} & \frac{1}{2} & \frac{1}{8} \end{array} \right]$$

→ Expand } 3-tap filter

Earlier times in 80s, working on image was studying of frequency stuff and signals. Still, some thing needs to be understood for CS in 2020s.

Fourier Transform

Idea: Every image/signal can be decomposed into set of things that describe things.

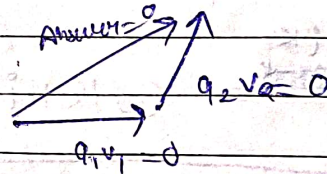
It is an important image processing tool which is used to decompose an image into its sine and cosine components.

⇒ Decomposing an image: basis sets (\hat{i}, \hat{j})

Suppose that $B = \{v_1, v_2, \dots, v_n\}$ is a finite subset of a vector space V over a field F (such as real or complex numbers). Then B is a basis if it satisfies:

- Linear independence

For all $a_1, \dots, a_n \in F$, if $a_1 v_1 + \dots + a_n v_n = 0$, then necessarily $a_1 = a_2 = \dots = a_n = 0$



- Spanning property

For every x in V it is possible to choose $a_1, \dots, a_n \in F$ s.t.

$$x = a_1 v_1 + \dots + a_n v_n = 0$$

* We need a good basis set, one is Fourier basis set.

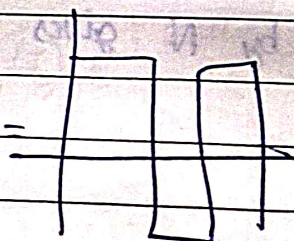
⇒ Fourier Theory:

Any periodic funcⁿ can be written as a sum of weighted sines and cosines of different frequencies.

$$g(t) = (A_1 \sin(\omega_1 t + \phi_1)) + A_2 \sin(\omega_2 t + \phi_2) + \dots$$

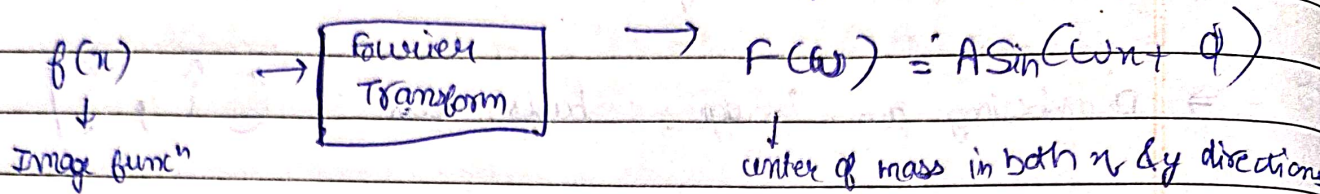
e.x.

Square wave



$$= A \sum_{k=1}^{\infty} \frac{1}{k} \sin(2\pi k t)$$

So, for decomposing our signal, we want to understand the frequency ω of our signal, phase is irrelevant as for CV we are not interested in reconstructing the image.



where $F(\omega) = R(\omega) + i I(\omega)$

$A = \pm \sqrt{R(\omega)^2 + I(\omega)^2}$

$\phi = \tan^{-1} \left(\frac{I(\omega)}{R(\omega)} \right)$

$F(\omega) = \int_{-\infty}^{\infty} f(x) e^{-i2\pi x \omega} dx$ (Fourier Transform)

$f(x) = \int_{-\infty}^{\infty} F(\omega) e^{-i2\pi x \omega} d\omega$ (Inverse Fourier Transform)

$e^{iK} = \cos K + i \sin K$

* $F(\omega)$ gives us the basis set.

How?

$\int_{-\infty}^{\infty} \sin(ax + \phi) \sin(bx + \phi) dx = 0$, if $a \neq b$



Independent

* $f(x) e^{-i2\pi x \omega}$ gives the winding machine exactly like 3Blue1Brown & integrating it (& dividing by N) gives centre of mass and ignoring divide by N gives

→ Discrete FT

$$F(k) = \frac{1}{N} \sum_{x=0}^{N-1} f(x) e^{-i \frac{2\pi kx}{N}}$$

where k is cycles per (Period of signal)

For 2-D

$$F(u, v) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) e^{-i2\pi(ux + vy)} dx dy \frac{1}{2}$$

$$F(k_x, k_y) = \frac{1}{N} \sum_{x=0}^{N-1} \sum_{y=0}^{N-1} f(x, y) e^{-i \frac{2\pi}{N} (k_x x + k_y y)}$$

Thus, we get the basis vector.

⇒ Fourier Transform & Convolution:

$$g = f * h$$

$$G(u) = \int_{-\infty}^{\infty} g(x) e^{-i2\pi ux} dx$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(\tau) h(x-\tau) e^{-i2\pi ux} d\tau dx$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} [f(\tau) e^{-i2\pi u\tau} d\tau] [h(x-\tau) e^{-i2\pi u(x-\tau)} dx]$$

$$= \int_{-\infty}^{\infty} [f(\tau) e^{-i2\pi u\tau} d\tau] \int_{-\infty}^{\infty} [h(x') e^{-i2\pi ux'} dx']$$

$$G(u) = F(u) H(u)$$

Convolution in spatial domain \Leftrightarrow multiplication in frequency domain

$$g = f * h \iff G = F \cdot H$$

$$g = f \cdot h \iff G = F * H$$

★ so, it's not un-common to calculate f & H and multiplying them & then using IFT to get g .

FFT is fast Fourier Transformation.

→ Smoothing & Blurring:

$$h(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{x^2}{2\sigma^2}}$$

$$H(u) = e^{-\frac{1}{2} (2\pi u)^2 \sigma^2} \quad \text{Also Gaussian.}$$

But, σ is multiplied instead of dividing.

$$G(u) = F(u) \cdot H(u) \quad \text{Thus, Attenuates Higher frequency.}$$

↓ IFT

ⓐ

⇒ Properties:

Linearity

$$c_1 f(x) + c_2 g(x)$$

$$c_1 F(u) + c_2 G(u)$$

Scaling

$$f(ax)$$

$$\frac{1}{|a|} F\left(\frac{u}{a}\right)$$

Differenti.

$$\frac{d^n f(x)}{dx^n}$$

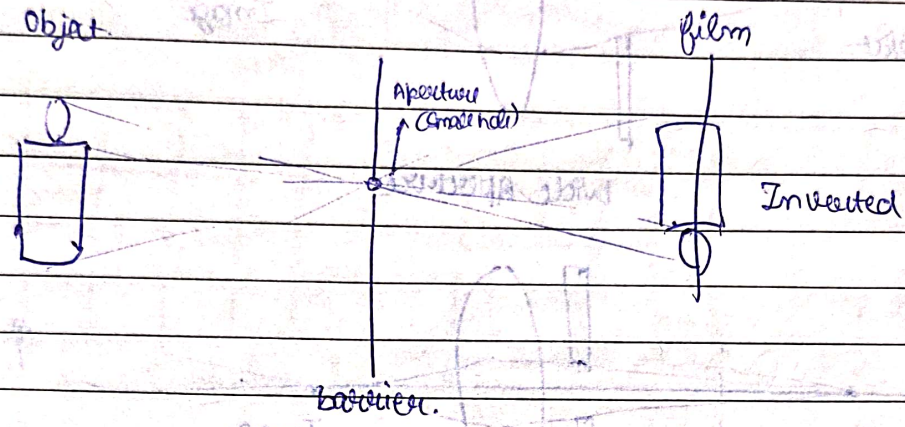
$$(i2\pi u)^n F(u)$$

3. Camera Models & Views:

Note, Image is a 2-D representation of 3D plane.

How to form an image?

a) Pinhole Perspective

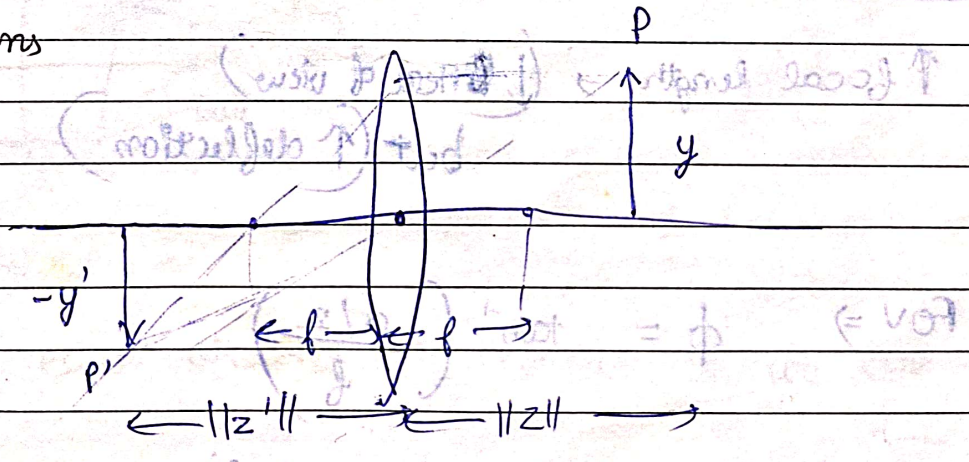


↓ Aperture Size → ↑ Crisp Image

But if hole is too small, diffraction takes place

↑ Aperture → ↓ Crisp Image

b) Thin Lens



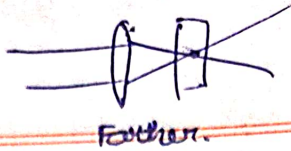
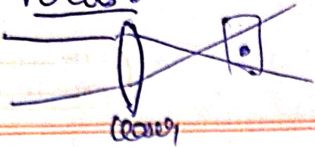
$$\frac{-y'}{y} = \frac{\|z'\|}{\|z\|} = \frac{\|z'\| - f}{f}$$

$$\Rightarrow \frac{\|z'\|}{\|z\|} = \frac{\|z'\| - f}{f}$$

$$\Rightarrow \frac{1}{f} = \frac{1}{u} + \frac{1}{v}$$

$u \rightarrow \|z'\|$
 $v \rightarrow \|z\|$

→ Focus:

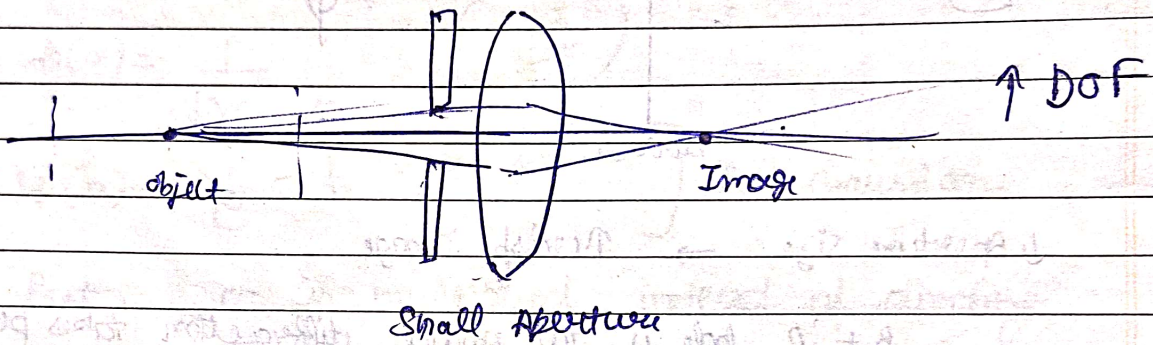
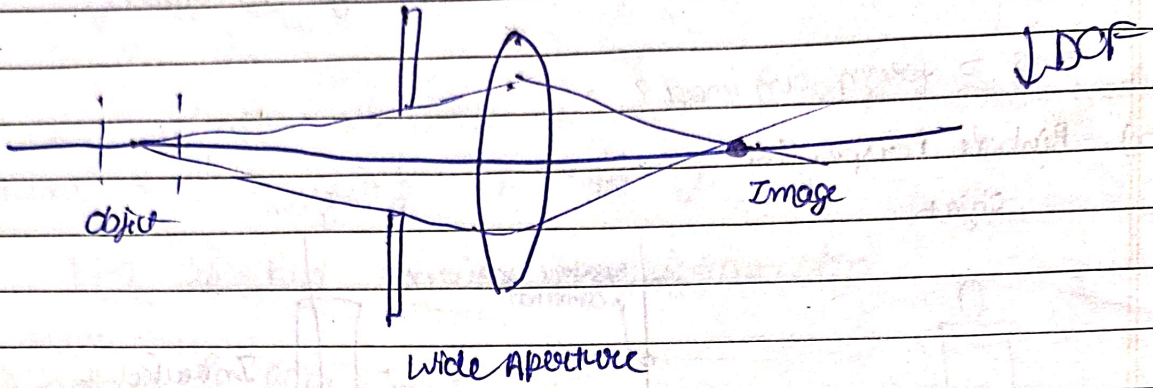


Date _____

Page _____



→ Depth of field:

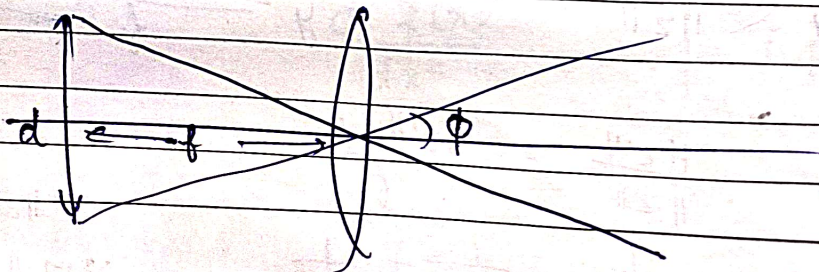


→ Field of view (zoom):

↑ focal length → (↓ Field of view)
but (↑ deflection)

$$FOV \Rightarrow \phi = \tan^{-1} \left(\frac{d/2}{f} \right)$$

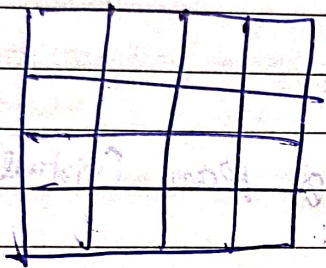
d → "retina" size or sensor size



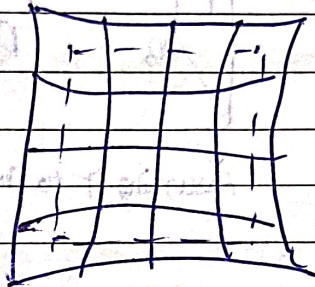
Zooming & moving are not the same

Problems of lens:
Lens are not perfectly thin

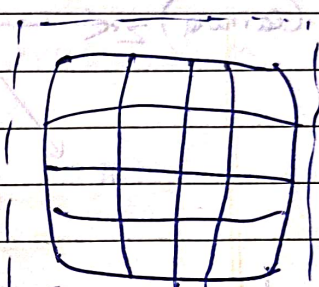
1. Geometric Distortion



No distortion



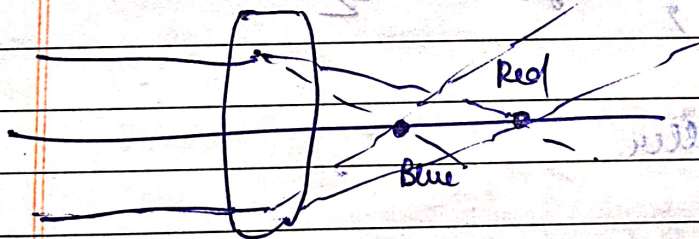
Pin cushion



Barrel

2. Radial distortion

3. Chromatic Aberration

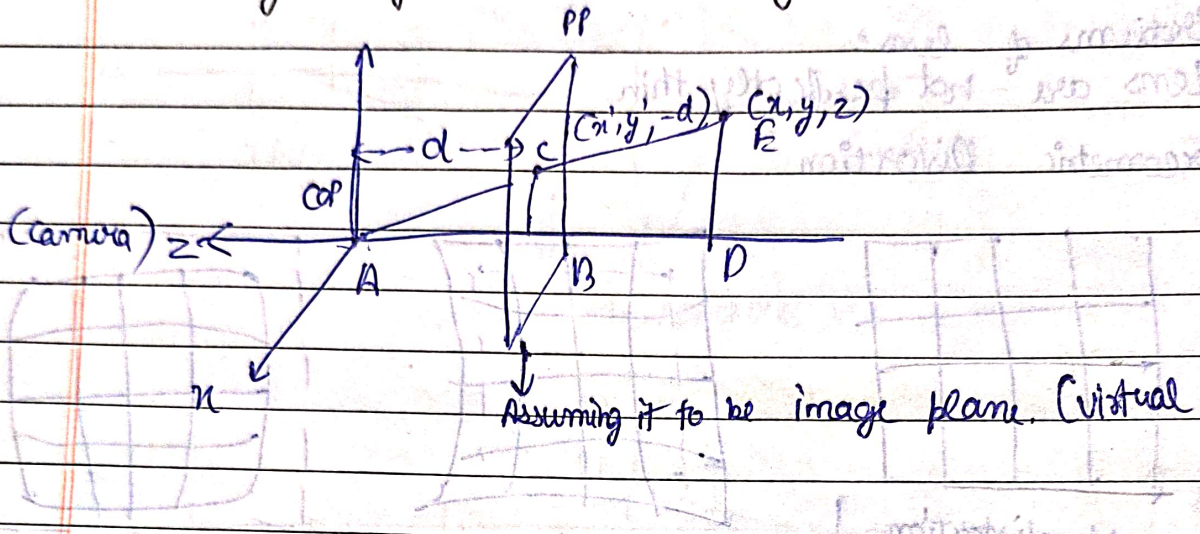


4. Vignetting

Dark edges around the corners in two lens system.

⇒ ~~Projecting~~ Perspective Imaging:

I. Modeling Projection - G-coordinate system



$\Delta ADE \sim \Delta ABC$

$$\frac{x'}{x} = \frac{y'}{y} = \frac{-d}{z}$$

$$\Rightarrow x' = -d \frac{x}{z}, \quad y' = -d \frac{y}{z}$$

∴ Distant objects are smaller.

Now $\frac{x}{z}$ & $\frac{y}{z}$ is not a linear transformation

So, add one more coordinate

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

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

Converting from homogeneous coordinate

For 2-D
 $x = \frac{x}{w}, \quad y = \frac{y}{w}$

For 3-D
 $x = \frac{x}{w}, \quad y = \frac{y}{w}, \quad z = \frac{z}{w}$

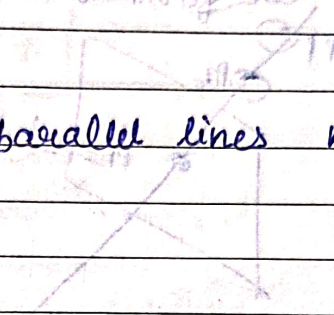
3D point to 2D image formula

$$\begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} fx \\ fy \\ z \end{bmatrix} = \boxed{x' = \frac{fx}{z}} \quad \boxed{y' = \frac{fy}{z}}$$

Here now, it is homogeneous

In perspective imaging, all parallel lines meet up at a vanishing point

e.g. Train tracks.



Mathematically,

A line in 3-space is represented as

$$x(t) = x_0 + at$$

$$y(t) = y_0 + bt$$

$$z(t) = z_0 + ct$$

$$x'(t) = \frac{fx}{z} = \frac{f(x_0 + at)}{z_0 + ct}$$

$$y'(t) = \frac{fy}{z} = \frac{f(y_0 + bt)}{z_0 + ct}$$

as $t \rightarrow \pm\infty$, for $c \neq 0$

$$\left. \begin{aligned} x'(t) &\rightarrow \frac{fa}{c} \\ y'(t) &\rightarrow \frac{fb}{c} \end{aligned} \right\}$$

So, ^{almost} all parallel (or not) lines converge to the same point.

Other Projection operators:

3-d point position

1. Perspective: $(x, y, z) \rightarrow \left(\frac{fx}{z}, \frac{fy}{z} \right)$
2. Weak Perspective: $(x, y, z) \rightarrow \left(\frac{fx}{z_0}, \frac{fy}{z_0} \right)$
3. Orthographic: $(x, y, z) \rightarrow (x, y)$



For State of the art algo.

Date _____

Page _____



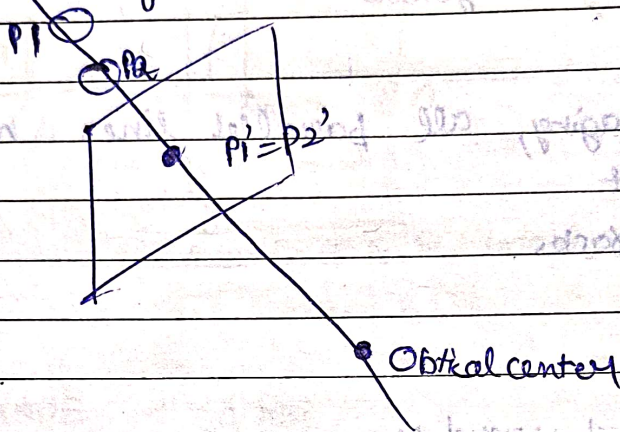
visit

vision.middlebury.edu/stereo

⇒ Stereo-Geometry (Geometry of Multiple views)

The depth of a scene point along with structure is not directly accessible in a single image.

B/c we lose information when converting an image from 3D to 2D



The trick to solving this would be to look at it through a different perspective.

eg. two eyes in animals, (vs)

* But, I can find depth with a single eye (am I non-human?)

Hey dumb-head, there are other factors such as

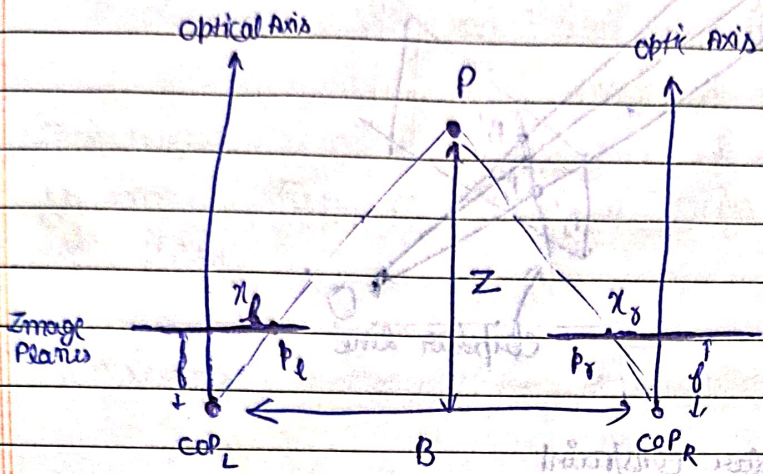
shading (shadow), texture, focus/de-focus, motion

Also, we have two eyes

Stereo:

- The image from one eye is ^{little} different than the image from the other eye.
- Think of shape from "motion" b/w ~~two~~ views.
- Infer 3-D shape of scene

• Basic Geometry



$$p_L = n_L \quad p_R = -n_R$$

Now (p_L, P, p_R) & (C_L, P, C_R) are similar

$$\text{So, } \frac{B - n_L + n_R}{B} = \frac{Z - f}{Z}$$

$$\Rightarrow Z = f \frac{B}{n_L - n_R}$$

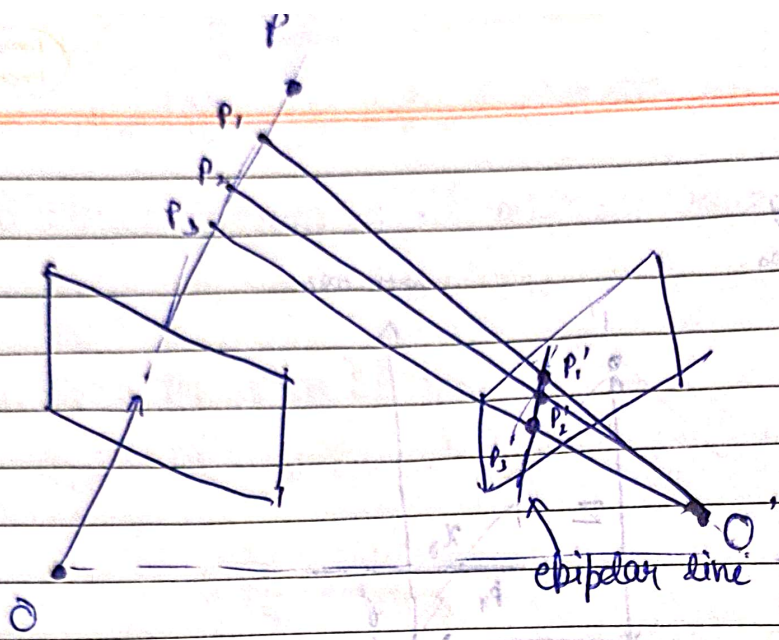
★ $n_L - n_R = \text{disparity}$

& $\text{disparity} \propto \frac{1}{\text{depth}}$

Now, we need to find disparity but we don't know either

n_L or n_R .

We will talk about epipolar geometry.



epipolar constraint

Given a calibrated stereo rig, the set of possible matches for the point p is constrained to lie on the associated epipolar line l' .

Now, did we solve the correspondence problem?

matching point in \mathcal{O}'

→ No, there are still a lot of points on the line

We need other "soft" constraints to help identify corresponding points

- Similarity
- Uniqueness
- Ordering
- Disparity gradient is limited.

1. Similarity

Assume image regions for the matches are similar in appearance.

For each pixel/window in the left image

- Compare with every pixel/window on the same epipolar line in right image
- Pick position with minimum match cost. (SSD, norm correlation)

Cross-correlation might not work good.

It depends a lot on window size.

2. Uniqueness

No more than one match in right image for each window in left image.

3. Ordering

abc in left image is abc in right image.

Won't always hold e.g. Consider transparent object.

e.g. Close one eye each time with both index fingers in front of nose and in a line.

* There are a lot of algorithms solving it, but there are still some challenges that stereo isn't possible to cover

- low-contrast; textureless image regions.
- Occlusions
- Violation of brightness constancy (e.g. specular reflections)
- larger Baseline.
- Camera Calibration errors.

Algo 1: Spline (DP) → Don't have stretching marks.

Algo 2: Full 2-D grid.
as Energy Minimization

$$E_{data} = \sum_i (w_1(i) - w_2(i + D(i)))^2$$

$$E_{smooth} = \sum_{\text{neighbours } i, j} p(D(i) - D(j))$$

$$\min. E = \alpha E_{data}(I_1, I_2, D) + \beta E_{smooth}(D)$$

World Co-ords



Camera Co-ords



film Co-ords



Pixel Co-ords

Rotation + Translation

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

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

Date _____

Page _____



⇒

Geometric Camera Calibration: We want to use the camera to tell us things about the world. So, we need the relationship between coordinates in the world & co-ordinates in the image.

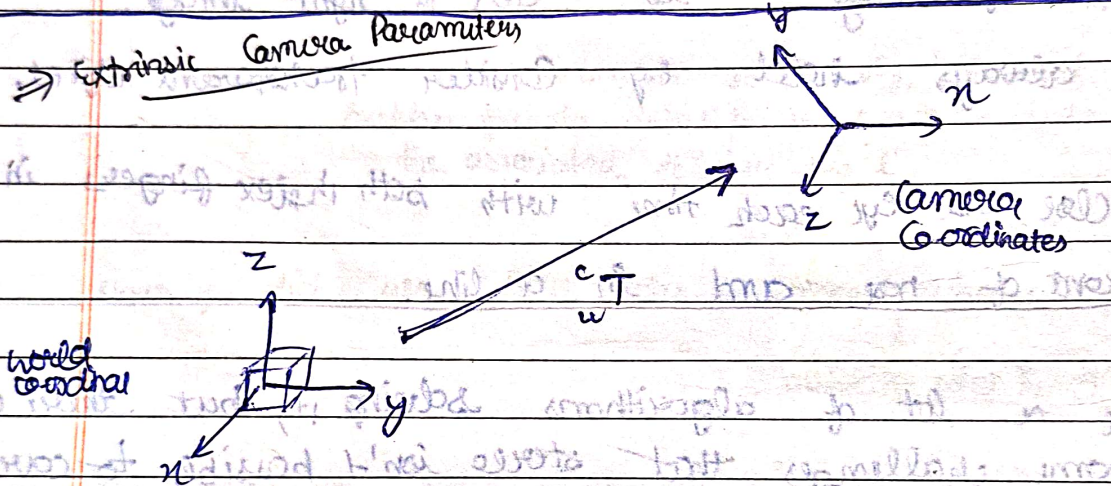
I.

Extrinsic Camera parameters: For arbitrary world co-ordinate system to the camera's 3D co-ordinate system. (Camera pos)

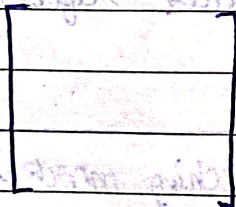
II.

Intrinsic Camera parameters: From the 3D co-ordinates in the camera frame to the 2D image plane via projection.

⇒ Extrinsic Camera Parameters



* A rigid body has 6 degrees of freedom (cube)



Translation

⇒ a point A in world coordinate is translated to point B in camera co-ordinate with a relation of

$${}^B P = {}^B (O_A) + {}^A P$$

↓
Origin A expressed in B frame

Commutative ✓

We can express the same eqⁿ using matrix multiplication, with the help of homogeneous co-ordinates

that means linear transformation

$$\begin{bmatrix} p_B \\ 1 \end{bmatrix} = \begin{bmatrix} I & B \\ O^T & 1 \end{bmatrix} \begin{bmatrix} p_A \\ 1 \end{bmatrix}$$

3x3 Identity 3x1

$$p_B = A p + B (O_A)$$

Rotation

$$\vec{OP} = \begin{bmatrix} i_A & j_A & k_A \end{bmatrix} \begin{bmatrix} A_x \\ A_y \\ A_z \end{bmatrix} = \begin{bmatrix} i_B & j_B & k_B \end{bmatrix} \begin{bmatrix} B_x \\ B_y \\ B_z \end{bmatrix}$$

$$B_p = \begin{pmatrix} B \\ A R \end{pmatrix} A_p$$

$$\begin{bmatrix} i_A \cdot i_B & j_A \cdot i_B & k_A \cdot i_B \\ i_A \cdot j_B & j_A \cdot j_B & k_A \cdot j_B \\ i_A \cdot k_B & j_A \cdot k_B & k_A \cdot k_B \end{bmatrix}$$

Also, by homogeneous co-ordinates.

$$\begin{bmatrix} B_p \\ 1 \end{bmatrix} = \begin{bmatrix} B & R & 0 \\ A & & \\ O^T & & 1 \end{bmatrix} \begin{bmatrix} A_p \\ 1 \end{bmatrix}$$

Commutative ~~X~~

Now, let's do a total rigid transformation,

→ A point 'A' in system is rotated and then translated (offset) in the A system to get point B in system

Using homogeneous co-ords

$$\begin{bmatrix} p_p \\ 1 \end{bmatrix} = \begin{bmatrix} A & B \\ R & O_n \\ A & 1 \\ O^T & 1 \end{bmatrix} \begin{bmatrix} A p \\ 1 \end{bmatrix} = \begin{bmatrix} B & A p \\ A & 1 \end{bmatrix}$$

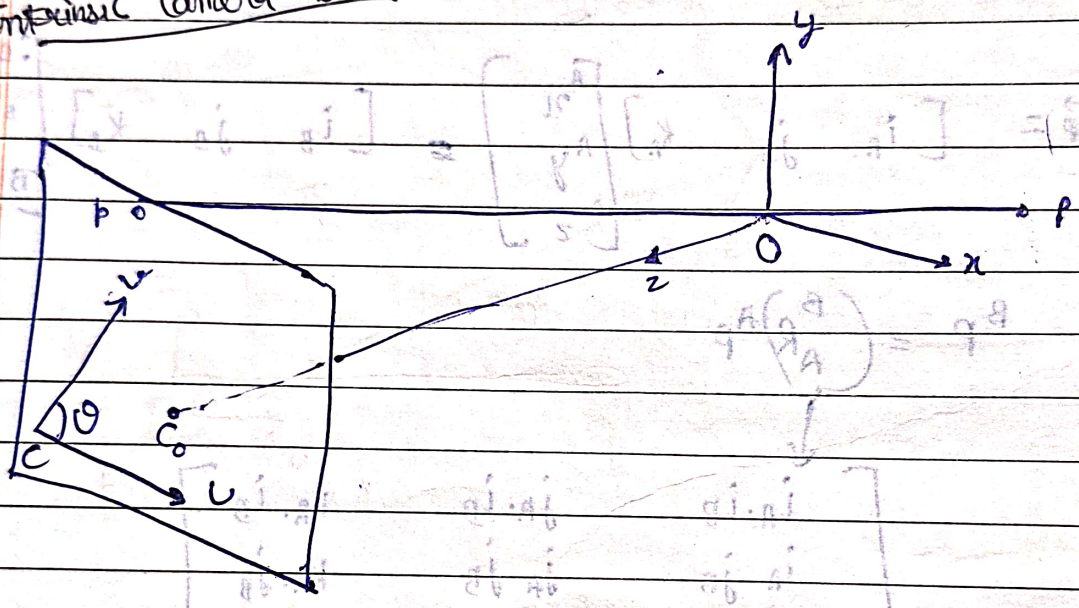
Point in camera frame

Extrinsic parameter matrix

Point in world frame

6 degrees of freedom

⇒ Intrinsic Camera Parameters



$$x' = u = \alpha \frac{x}{z} + u_0$$

$$y' = v = \beta \frac{y}{z} + v_0$$

B/c α is not in mm α is there
 now $\alpha \neq \beta$ so, we have two degrees of freedom.

origin of camera co-ordinate system is at corner 'c' of setting and not at it's center so u_0 & v_0
 4 degrees of freedom.

(Principal points)

Also, the camera coordinate system may be skewed due to some manufacturing error, so angle θ b/w two image axes is not equal to 90° .

$$\begin{aligned} u &= \alpha \frac{x}{z} - \alpha \cot \theta \frac{y}{z} + u_0 \\ v &= \frac{\beta}{\sin \theta} \frac{y}{z} + v_0 \end{aligned}$$

So, 5 degrees of freedom.

$$\begin{bmatrix} z \cdot u \\ z \cdot v \\ z \end{bmatrix} = \begin{bmatrix} \alpha & -\alpha \cot \theta & u_0 & 0 \\ 0 & \beta \cot \theta & v_0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

$$\vec{p}' = K \vec{c}_p$$

$$\Rightarrow \text{Actual point} = \frac{\vec{p}'}{z}$$

→ Combining extrinsic & intrinsic calibration parameters

$$\vec{p}' = K \begin{pmatrix} c & R \\ w & t \end{pmatrix} \vec{w}_p$$

$$\vec{p}' = M \vec{w}_p$$

↓ pixel coords
 ↑ World coords

where

$$M_{(3 \times 4)} = \begin{bmatrix} f & s & x_c' & 0 \\ 0 & af & y_c' & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} R_{3 \times 3} & 0_{3 \times 1} \\ 0_{1 \times 3} & 1 \end{bmatrix} \begin{bmatrix} I_{3 \times 3} & T_{3 \times 1} \\ 0_{1 \times 3} & 1 \end{bmatrix}$$

intrinsic Projection Rotation Translation

⇒ CAMERA CALIBRATION:

I. Calibration using known points (Direct linear transformation)

Place a known object in the scene.

- Identify correspondences between image and scene
- Compute mapping from scene to image.

~~Identify correspondences~~

⇒ Now, we need to find all 11 degrees of freedom. It is done using SVD tricks.

- Find the m that minimizes $\|Am\|$ subject to $\|m\|=1$
- Let $A = UDV^T$ (SVD)

$\frac{1}{\|m\|}$
original
the

- Therefore minimizing $\|UDV^T m\|$.

- But, $\|UDV^T m\| = \|DV^T m\|$ & $\|m\| = \|V^T m\|$

U & V are unit vectors

- Thus, minimize $\|DV^T m\|$ subject to $\|V^T m\|=1$.

Let $y = V^T m$

So, $\|Dy\|$ is minimized s.t. $\|y\|=1$

Also, D is diagonal with decreasing values.

So $\|Dy\|$ minimum is when $y = (0, 0, \dots, 0, 1)^T$

Since $y = V^T m$, $m = Vy$ as V is orthogonal ($V^{-1} = V^T$)

This gives last column in V .

* & $A^T A = V^T D^2 V$ (eigen decomposition)

The final solⁿ is:

Given $Am=0$, find the eigenvector of $A^T A$ with smallest eigenvalue, that's m .

⇒ After finding M , we should be able to find things like the camera center from M

Two ways: Pure Way
Easy Way

$$M = [Q | m_4]$$

$$Q = 3 \times 3$$

$$m_4 = 3 \times 1$$

$$C = -Q^{-1} m_4$$

II. Multi-plane Calibration.

Advantages

- Requires a plane only (checker-board)
- Don't need to know positions/orientations
- Direct code available in open CV

Zhengyu Zhang (Good of Calibration)



⇒ Multiple Views

⇒ Image to Image Mappings

$$\text{I} \quad \begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} \approx \begin{bmatrix} \alpha x' \\ \beta y' \\ 1 \end{bmatrix}$$

6 degree of freedom

Affine transformation

	D.O.F.
↳ Translation	2
↳ Rotation	1
↳ Scale	1
↳ Skew	2

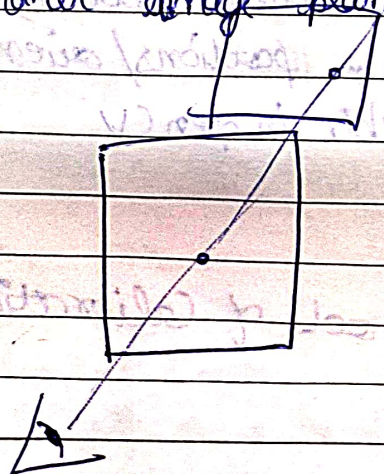
$$\text{II} \quad \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} x \\ y \\ w \end{bmatrix} = \begin{bmatrix} w'x' \\ w'y' \\ w'1 \end{bmatrix} \approx \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} \quad (\text{Divide by } w')$$

Homography

8 degree of freedom (This transformation maps for any linear ~~between~~ between any two planes)

⇒ Homographies and mosaics

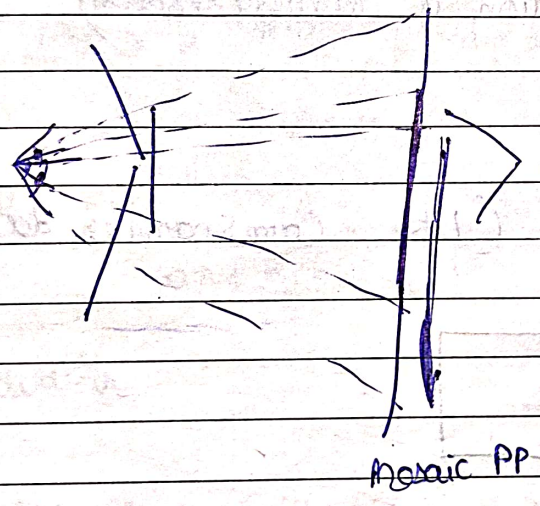
Rather than thinking of this as a 3D reprojection on 2D plane, we surely can think of it as a 2D image warp from one image plane to another image plane.



This gives an idea of panoramas.

→ How to stitch together a panorama.

- Take a sequence of images from the same position.
→ rotate the camera about its optical center.
- Compute transformation b/w second image plane & first.
- Transform the second image to overlap with the first.
- Blend the two together to create a mosaic.
(Repeat, if more than two)



The mosaic has a natural interpretation in 3D: Overlap

- The images are reprojected onto a common plane
- The mosaic is formed on this plane.

(works for views upto 180°, Beyond that need to use sequence that "bends the rays" or map onto a different surface, say, a cylinder)

→ Now, how to solve homography?

- Since, homography is transformation b/w any two planes,

$$P' = HP$$

$$\begin{bmatrix} wx' \\ wy' \\ w \end{bmatrix} = \begin{bmatrix} a & b & e \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

(8 degrees)

- Requires at-least 4 points for 8 eqs.

Non-homogeneous way
 $Ah = b$

(h is vector of 8 degree of freedom)

Solve for h by $\min \|Ah - b\|^2$ using least-squares.

OR

Homogeneous way

- Solve using SVD

→ Applications of Homography (Not mosaic dummy-dumb)

i) Video motion (Computation of camera motion)

ii) Image Rectification

Rectifying slanted views (Like CamScanner does)



How do we do this?

⇒ Projective Geometry

In case of homogeneous co-ordinates

- Line joining two points p_1 & p_2 is given by

$$p_1 = [x_1, y_1, 1]$$

$$p_2 = [x_2, y_2, 1]$$

$$l = p_1 \times p_2$$

- Similarly, intersection of two lines

$$l_1 = [a_1, b_1, c_1]$$

$$l_2 = [a_2, b_2, c_2]$$

$$p_2 = l_1 \times l_2$$

Before that,

Projective line

(in image plane)

A line is a plane of rays through define by the normal line

$$l = (a, b, c)$$

All rays (x, y, z) satisfying:

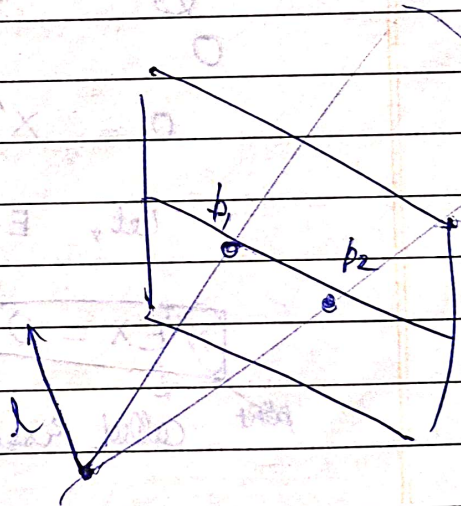
$$ax + by + cz = 0$$

Point & line duality

- A line l is a homogeneous 3-vector

l is \perp to p_1 & p_2

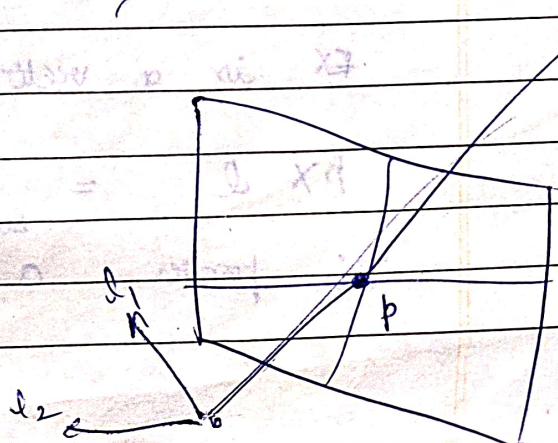
$$l = p_1 \times p_2$$



- A point p is \perp to l_1 & l_2

$$p = l_1 \times l_2$$

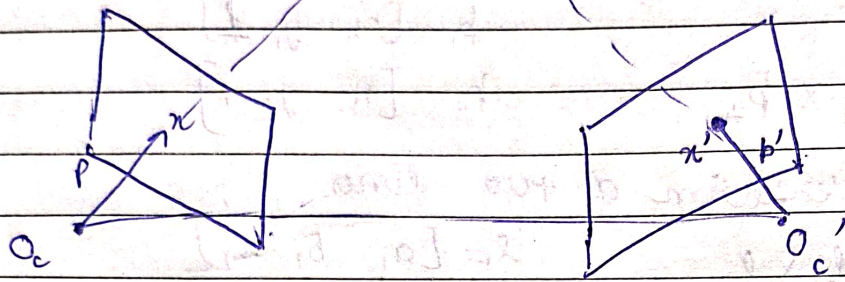
(Go to top now)





Essential Matrix

Consider again the stereo problem,
a world point



In case of calibrated cameras

$$x'_c = R x_c + T$$

$$T x x' = T x R x + T x T$$

$$T x x' = T x R x$$

$$x' \cdot (T x / x') = x' \cdot (T x R x)$$

0

$$0 = x' \cdot (T x R x)$$

Let, $E = T x R$

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

point called essential matrix

$$x E x'^T = 0$$

$E x$ is a vector l (line)

$P x l$ = line

This proves a point, lies on a epipolar line

⇒ Fundamental Matrix

Motivation:

- estimate epipolar geometry from a set of point correspondences between two uncalibrated cameras.

We know,
$$P_{im} = K_{int} \underbrace{\Phi_{ext}}_{P_c} P_w$$

$$\Rightarrow P_{im} = K_{int} P_c$$

$$\Rightarrow \textcircled{P_c} = K_{int}^{-1} P_{im} \quad (K_{int} \text{ is invertible})$$

Camera → Image

this gives the ray in space (as it is homogenous)

and for two cameras

$$P_{c, \text{left}} = K_{int, \text{left}}^{-1} P_{im, \text{left}}$$

$$P_{c, \text{right}} = K_{int, \text{right}}^{-1} P_{im, \text{right}}$$

$$P_{c, \text{right}}^T E P_{c, \text{left}} = 0 \quad (\text{look at it})$$

$$(K_{int, \text{right}}^{-1} P_{im, \text{right}})^T E (K_{int, \text{left}}^{-1} P_{im, \text{left}}) = 0$$

$$P_{im, \text{right}}^T \underbrace{(K_{int, \text{right}}^{-1} E K_{int, \text{left}}^{-1})}_F P_{im, \text{left}} = 0$$

$$P_{im, \text{right}}^T F P_{im, \text{left}} = 0$$

⇒ $l = F p'$ is the epipolar line in the b image associated with b' & vice versa ($l' = F^T p$)

→ Properties of Fundamental Matrix.

- Epipolar found by $Fp' = 0$ & $F^T p = 0$

↓
a point which exist on every epipolar line

* In parallel images planes, epipolar are at infinity.

- F is singular (mapping from homogeneous 2-D point to 1-D family so rank 2)

- If we estimate fund. matrix from correspondences in pixel coords, can re-construct epipolar geometry without intrinsic or extrinsic parameters.

→ How to find it? (Notes: FM gives relation btw Any two views That gives an idea of video motion)

- Each point generates one constraint or single eqⁿ on F .

$$p^T F p' = 0$$

Not famous.

Requires 8 point atleast to solve for F (9 variables)
divide 1 variable
↓
8 variables left
Now, do OLS or orthogonal least squares.

more used [SVD approach (very minimal errors)

$F =$ Find eigen vector of $A^T A$ with minimum eigen value or OLS

Do another SVD (Since rank of F is 2)

$$F = U D V^T$$

$\hat{D} =$ set the last singular value to zero.

→ Applications of Stereo:

- Stereo Rectification
- Photo-Synth.

→ Summary:

- For 2-views, there is a geometric relationship between rays in one view to rays to the other → epipolar geometry.
- These relationships can be captured algebraically as well
 - Calibrated: Essential Matrix
 - Uncalibrated: Fund. Matrix
- This relation can be estimated from point correspondences

→ Another problem?

How to find these point correspondences?
manually.....

Not

We will use features to find these.

$$\begin{bmatrix} u' & v' & 1 \end{bmatrix} F \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = 0$$

$$\rightarrow \begin{bmatrix} u'_1 u_1 & u'_1 v_1 & u'_1 & v'_1 u_1 & v'_1 v_1 & v'_1 & u'_1 - v'_1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ u'_N u_N & u'_N v_N & u'_N & v'_N u_N & v'_N v_N & v'_N & u'_N - v'_N & 1 \end{bmatrix} F = 0$$

N points

$$A n = 0$$

$$P = U \hat{D} V^T$$

Features & Matching

⇒ Introduction to "features":

To get a transformation of one image to another image, we need corresponding points, and we do that using "LOCAL FEATURES".

Goal: Find points that are:

→ Found in other images

→ Found precisely

→ Found reliably.

→ Characteristics of Good Features:

- Precision / Repeatability

The same feature can be found in several images despite geometric & photometric transformations.

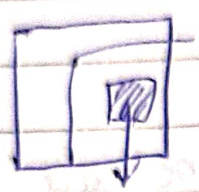
- Saliency / Matchability

Each feature, has a distinctive property.

- Compactness / Efficiency

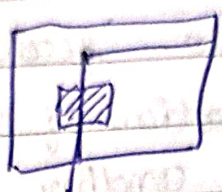
Many fewer features than image pixels.

⇒ Finding Corners:



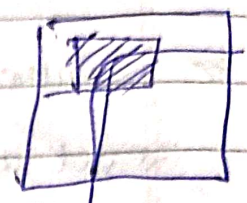
Flat

No change in all directions



Edge

No change along the edge direction



corner

significant change in all direction with small shift

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

Now, in this eqⁿ, we want to look for small shifts
↓ (Derivative)

second-order Taylor expression of $E(u,v)$ about $(0,0)$

For 1-D

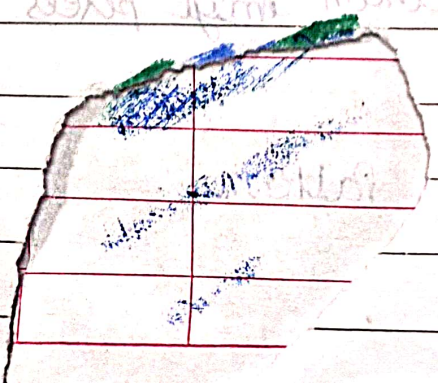
$$f(s) \approx f(0) + s \frac{df(0)}{dx} + \frac{1}{2} s^2 \frac{d^2f(0)}{dx^2}$$

For 2-D

$$E(u,v) \approx E(0,0) + [u \ v] \begin{bmatrix} E_u(0,0) \\ E_v(0,0) \end{bmatrix} + \frac{1}{2} [u \ v]$$

$$\times \begin{bmatrix} E_{uu}(0,0) & E_{uv}(0,0) \\ E_{uv}(0,0) & E_{vv}(0,0) \end{bmatrix}$$

$$\times \begin{bmatrix} u \\ v \end{bmatrix}$$



Now,

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

$$E_{UU}(u, v) = \sum_{x,y} 2w(x,y) I_x(x+u, y+v) I_x(x+u, y+v)$$

$$+ \sum_{x,y} 2w(x,y) [I_x(x+u, y+v) - I_x(x,y)] I_{xx}(x+u, y+v)$$

$$E_{UV}(u, v) = \quad \quad \quad I_y \quad \quad \quad I_x$$

Putting all in eqⁿ, we get

$$E(u, v) = \frac{1}{2} [u \ v] \begin{bmatrix} \sum_{x,y} w(x,y) I_x^2(x,y) & 0 \\ I_x I_y & \sum_{x,y} w(x,y) I_x(x,y) I_y(x,y) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

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

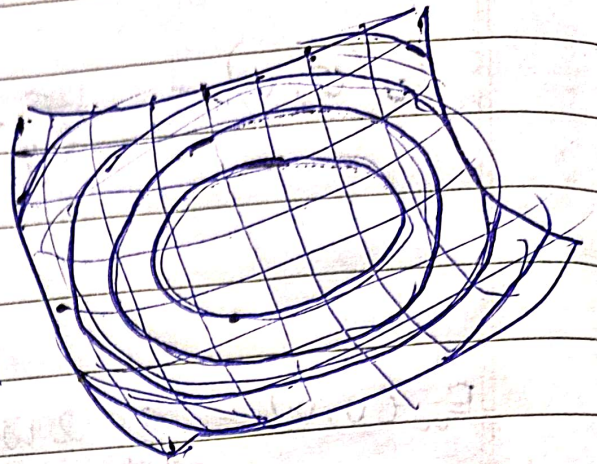
$$\Rightarrow E(u, v) \approx [u \ v] M \begin{bmatrix} u \\ v \end{bmatrix}$$

Interpreting it, (consider a slice of $E(u, v)$)

$$E(u, v) = \sum I_x^2 u^2 + 2 \sum I_x I_y uv + \sum I_y^2 v^2 = K$$

The eqⁿ is of ellipse.

- * Treat gradient vectors as a set of (dx, dy) points with a center of mass defined as being at $(0, 0)$
- * Fit an ellipse to ^{that} set of points via scatter matrix.



Now again consider the \bar{m} matrix (second moment)

$$\# \quad M = \sum_{n,y} w(x,y) \begin{bmatrix} I_n^2 & I_n I_y \\ I_n I_y & I_y^2 \end{bmatrix}$$

Case 1a Gradients are either horizontal or vertical but not slanted

$$I_n I_y = 0$$
$$M = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

$$\left. \begin{array}{l} \lambda_1 \\ \lambda_2 \end{array} \right\} \begin{array}{l} a \text{ \& } b \\ \text{for ellipse} \end{array}$$

Case 1b Gradient is only horizontal: or only vertical
either $\lambda_1 = 0$ or $\lambda_2 = 0$
(Not a good corner)

Notes: λ_1 & λ_2 determine the corner based on the property of diagonalization of M matrix
How? and concluding them as eigen values.
Watch video again.

λ_1 & λ_2 are eigen-values of M ,

Case I: λ_1 & λ_2 are too small

Flat region

Case II: $\lambda_1 \gg \lambda_2$ or $\lambda_2 \gg \lambda_1$

Edge

Case III: λ_1 & λ_2 are large, $\lambda_1 \sim \lambda_2$
Corner.

But, Harris algorithm makes use of "response function"

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

$$R = \lambda_1 \lambda_2 - \alpha (\lambda_1 + \lambda_2)^2$$

and α : constant (0.04 to 0.06)

So, $|R|$ small \rightarrow Flat

$R < 0$ \rightarrow Edge

$R \gg 0$ \rightarrow Corner

\rightarrow Harris detector: Algorithm

1. Compute gaussian derivatives at each pixel.
2. Compute second moment matrix M in a gaussian window around each pixel.
3. Compute corner response funcⁿ R .
4. Threshold R .
5. Find local maxima of response function (non-maximum suppression)

Thus, we get good features on both images.

• Properties of Harris detector:

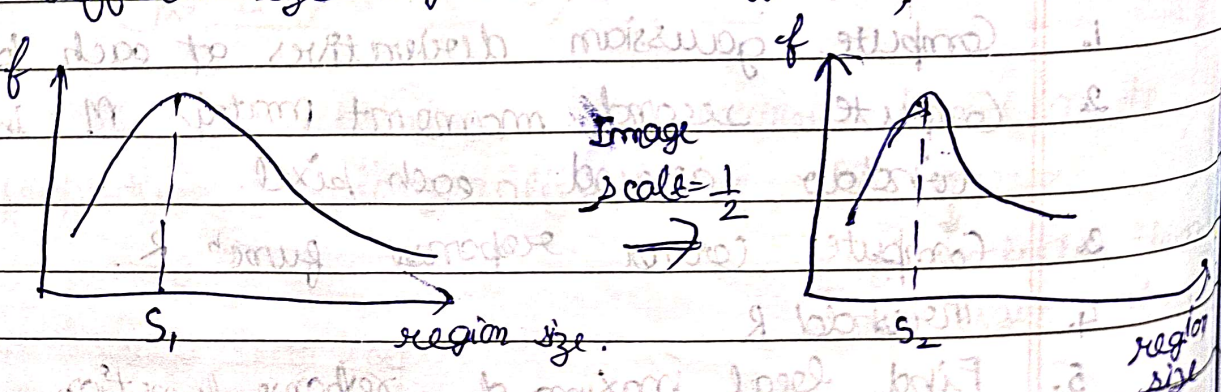
1. Invariant to Rotation.
2. Intensity
 - Invariant to Addition
 - Threshold Issue in Multiplication but invariant.
3. Not invariant to scale.

How to do scale invariance?

Solⁿ: Design a funcⁿ on the region, which is "scale invariant" - not affected by the size but will be the same for "corresponding regions", even if they are at different scales.

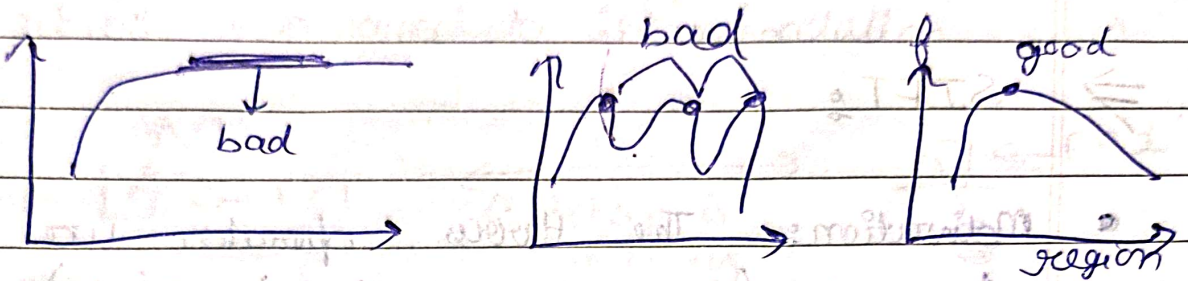
One Method:

At a point, compute the scale invariant funcⁿ over different size neighborhood (different scales)



we get scale size S_1 & S_2

Another problem, we need a good funcⁿ for scale invariant detection, that has one stable sharp peak.



* A good funcⁿ would be a one which responds to contrast.

$$L = \sigma^2 (G_{xx}(x, y, \sigma) + G_{yy}(x, y, \sigma))$$

↓
Laplacian of Gaussian - LoG

It's a little costly to find.

we instead use DOG, that is essentially the same

$$DOG = G(x, y, K\sigma) - G(x, y, \sigma)$$

↓
Difference of Gaussian

Both
Kernels are
invariant
to scale
& rotation

Now, we get interest points but we still need to match them ^{distinctively}, that is the work of descriptor.

⇒ SIFT:

- Motivation: The Harris operator was not invariant to scale (but we fixed that in 2004) and correlation was not invariant to rotation
↓
(normal Brute force matching)

- Develop

Done / Detector → invariant to scale & rotation

Descriptor → robust to variations in viewing conditions.

- Idea of SIFT

Match enough SIFT features that are invariant to scale, transformations, photometry etc. to do computer vision recognition.

- Overall SIFT procedure

1. Scale-space extrema detection

2. Keypoint localization

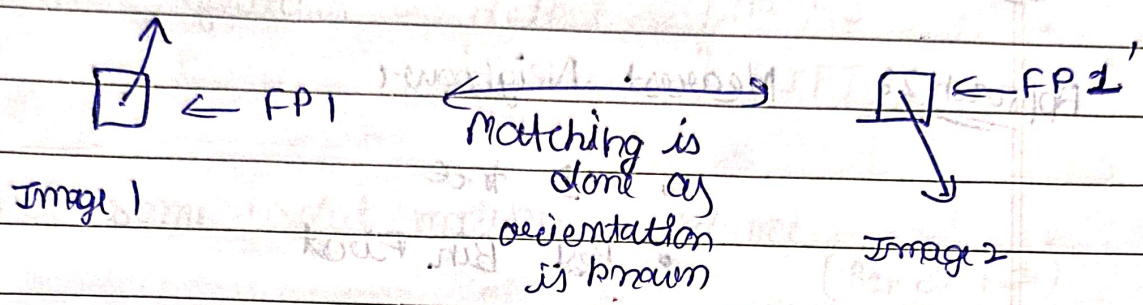
3. Orientation Assignment

4. Keypoint Description

} Or we Harris-Laplace or other method

→ Orientation Assignment

An orientation is assigned to each keypoint to achieve invariance to image rotation.



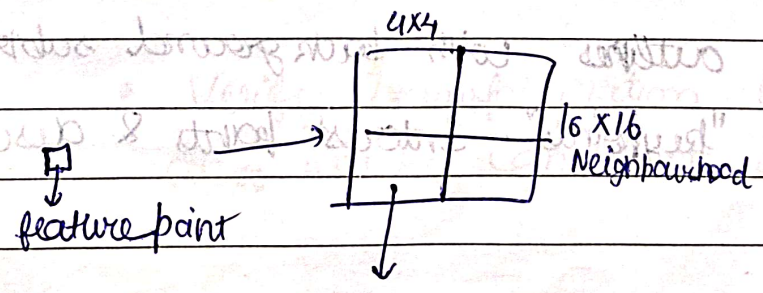
→ Keypoint Descriptors

Next is compute a descriptor for the local image region about each keypoint that is:

- Highly distinctive
- As invariant as possible

But first Normalization...

- Rotate to standard orientation
- Scale the window size



Histogram with 8 bins = 128 bin values called as **SIFT** feature vector

Now, we know how to describe them but
 how about matching? (we still not done 😞)

Approach 1: n^2 Brute force

Approach 2: Nearest Neighbour

- k-D tree
- Best Bin First

Approach 3: Wavelet-Based Matching (Hashing) ←

↓
 • Compute a short (3-vector) descriptor from the neighbourhood

- Quantize each value into 10 (overlapping) bins
 10^3 entries.

⇒ 3D Object Recognition.

Train

1. Extract outlines with background subtraction
2. Compute "keypoints" - interest points & descriptors.

Test

1. Find possible matches
2. Search for consistent solution
 - such as affine.

⇒ Feature-Based Alignment :

→ Overall strategy

- Compute features - detect } Harris
describe } SIFT
- Find some-useful matches - kd-tree
BBF (Best Bin First)
- Compute & apply the best transformation - Affine
Homography

Algorithm:

1. Extract features
2. Compute putative matches → "closest descriptor"
Not known how to do
3. Loop until happy :
 - Hypothesize transformation T from some matches
 - Verify Transformation (search for other matches consistent with T) - mark best

Here step 2 isn't encountered yet, so tackling it now.

* We know exhaustive search, Hashing & Nearest Neighbor techniques gives the best match, how do we know it's a good one?

How do we know a good match is found?

Problem!

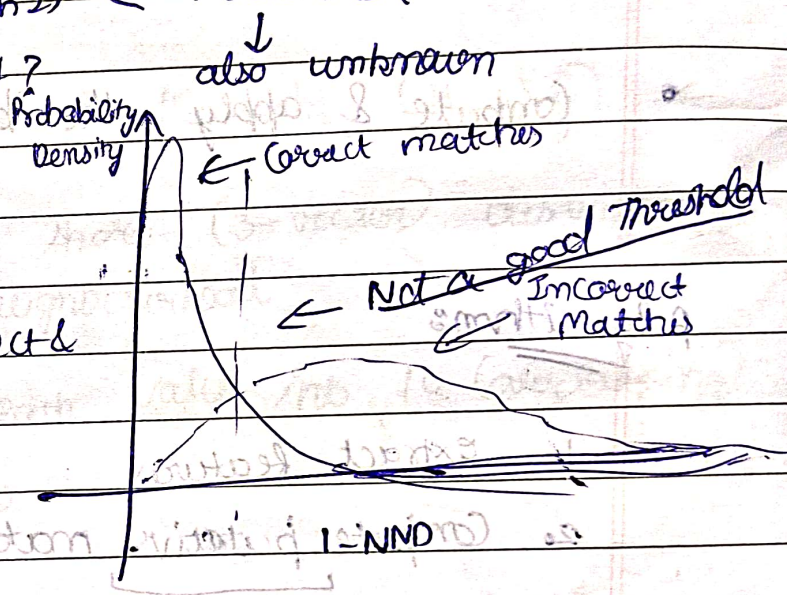
Feature Space Outlier Rejection:

Let's not match all features, but only those that have "similar enough" matches?

How can we do it?

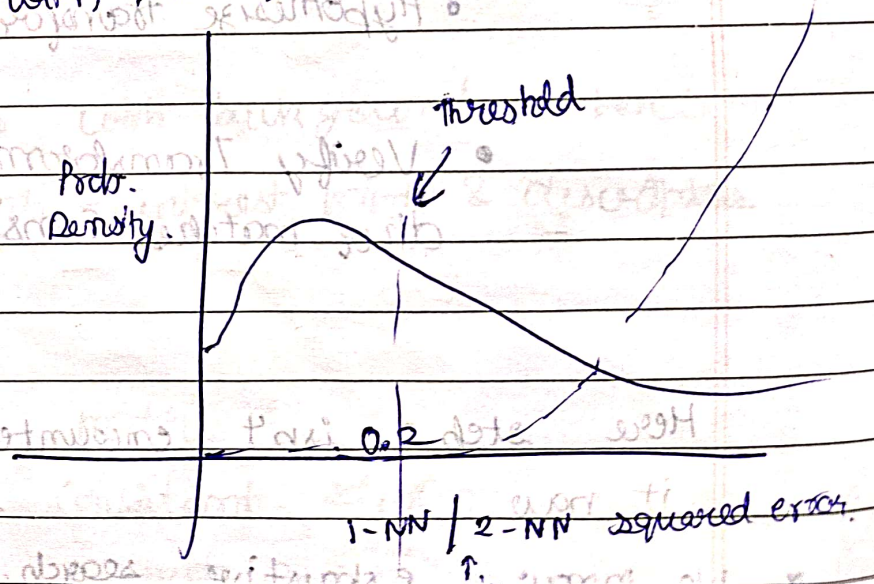
SSD (patch1, patch2) < threshold

But, how to threshold it?



A point's Best matches are sometimes correct & sometimes not.

So, we use Lowe's method to use 2 NND and compare it with 1-NND.



combine

Problem 2

even when pick best match, still lots (& lots) of wrong matches - "outliers", because we did the distinctive property thing.

essentially, we need to fit a model s.t. the effect of outliers is minimized.

⇒ RANSAC (Random Sample Consensus)

- Randomly pick some points to define your line (model). Repeat enough times until you find a good line (model) - one with many inliers.
- It somehow, copes with a large proportion of outliers.
- Every general model has a minimal set,
for line it is 2
for image transformations, 8 points.

• Translation: One pair of matched points

• Homography (for plane): 4

• Fundamental Matrix: 8 point pairs (in reality, 7)

$$(pp \cdot O = 0 \quad p \cdot O)$$

$$M \cdot H = 0$$

Algorithm

1. Randomly select s points (or point pairs) to form a sample.
2. Instantiate the model.
3. Get consensus set C_i - the pairs within error bounds (distance threshold) of the model.
4. If $|C_i| > T$, terminate & return model.
5. Repeat for N trials, return model with max. $|C_i|$.

Parameters

1. s - Initial no. of points \rightarrow known
2. Distance threshold t , \rightarrow unknown.
 - Choose t so prob. for inlier is high 95%.
 - 0 mean gaussian noise with std. dev. σ
$$t^2 = 3.84 \sigma^2$$
3. Number of trials N .
 - Choose N so that, with probability p , at least one random sample set is free from outliers (e.g. $p = 0.99$)
 - Need to set N based upon the outlier ratio ϵ

→ Calculating N

1. $n \rightarrow$ no. of points to compute solution

2. $p \rightarrow$ prob. of success.

3. $e \rightarrow$ proportion of outliers, so % inliers = $(1-e)$

4. $P(\text{sample set with all inliers}) = (1-e)^n$

5. $P(\text{sample set will have at-least one outlier})$

$$= 1 - (1-e)^n$$

6. $P(\text{all } N \text{ samples have outlier}) = (1 - (1-e)^n)^N$

7. We want $P(\text{all } N \text{ samples have outlier}) < (1-p)$

$$\Rightarrow (1 - (1-e)^n)^N < (1-p)$$

$$\Rightarrow \boxed{N > \frac{\log(1-p)}{\log(1 - (1-e)^n)}}$$

→ lastly, do re-compute least-squares fit estimate on all of the inliers.

Adaptive algorithm

Set $N = \infty$, $sample_cnt = 0$, $e = 1.0$

while $N > sample_cnt$

Choose s & count the no. of indices

st $e_0 = 1 - \frac{\text{no. of indices}}{\text{No. of points}}$

if $e_0 < e$, set $e = e_0$ & recompute

N from e :

$$N = \frac{\log(1 - \epsilon)}{\log(1 - (1 - \epsilon)^2)}$$

$sample_cnt++$

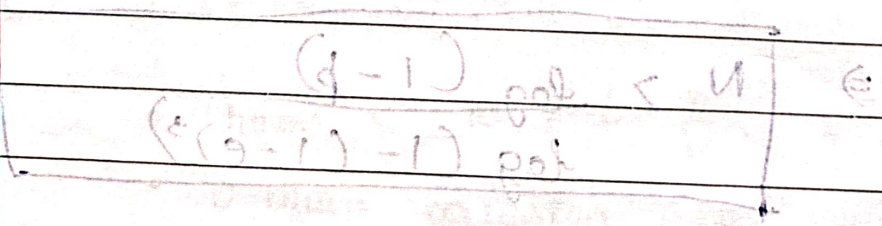


PHOTO METRY

factors affecting image:

1. Lights, Surfaces & Shadows
2. Reflections
3. Refractions
4. Inter-reflections
5. Scattering.

→ Surface Appearance

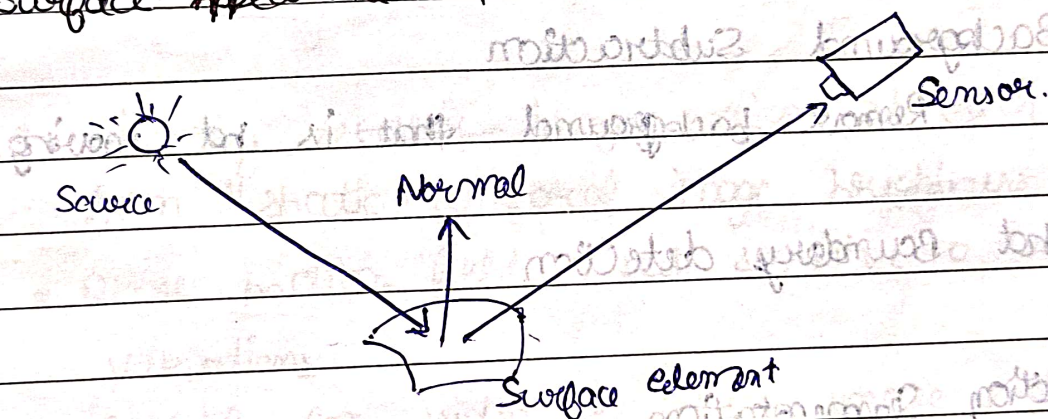


Image Intensity = f (Color, surface reflectance, illumination)

- Surface reflection depends on both the viewing & illumination directions.

⇒ Introduction to Motion

- A video is a sequence of frames captured over time - usually quickly.

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

\downarrow \downarrow
 space time

- Motion Applications:

→ Background Subtraction

Remove background that is not moving.

→ Shot Boundary detection

→ Motion Segmentation

Segment the video into multiple coherently moving objects.

Goal: Manufacture information in motion like brain does.

- Estimating 3D structure
- Learning dynamic models
- Recognizing events & activities
- Improving video quality.

→ Motion estimation techniques :

I. Feature-based Methods (Done ✓)

- Extract visual features & track them over multiple frames.
- Sparse motion fields, but more robust tracking
- Suitable when image motion is large.

II. Direct, dense methods

- Directly recover image motion at each pixel from "spatio-temporal image brightness variation."
- Dense motion fields, but sensitive to appearance variations
- Suitable for video & when image motion is small.

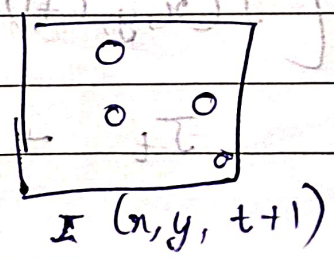
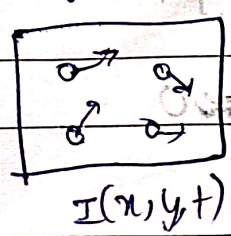
State of the art → Combination of both above.

→ Dense Flow : Brightness Constraint

i) Optic Flow

It is the apparent motion of objects or surfaces.

Problem Defⁿ:



How to estimate pixel motion from image 1 to 2.

⇒ solve pixel co-correspondence problem

Given a pixel in $I(x, y, t)$, look for nearby pixels of the same color in $I(x, y, t+1)$

This is the optic flow problem.

Assumptions

• color constancy $I(x, y, t) = I(x+u, y+v, t+1)$

For grayscale, this is brightness constancy

• small motion. u, v are less than 1 px, or smooth

$$I(x+u, y+v) = I(x, y) + \frac{\partial I}{\partial x} u + \frac{\partial I}{\partial y} v + \dots$$

↑
higher order

Combining both eqn.

$$I(x+u, y+v, t+1) - I(x, y, t) = 0$$

$$I(x, y, t+1) + I_x u + I_y v - I(x, y, t) \approx 0$$

(Assume $t \sim$ small)

$$\Rightarrow [I(x, y, t+1) - I(x, y, t)] + I_x u + I_y v \approx 0$$

$$\Rightarrow I_t + I_x u + I_y v \approx 0$$

$$0 \approx I_t + \nabla I \cdot \langle u, v \rangle$$

$$\nabla u, v \rightarrow 0$$

this becomes exact eqⁿ.

$$I_t + \nabla I \cdot \langle u, v \rangle = 0$$

Brightness constancy constraint eqⁿ.

* 2 unknowns, but 1 eqⁿ
(u, v)

similar to

$$au + bv + c = 0$$

I_x, I_y, I_t

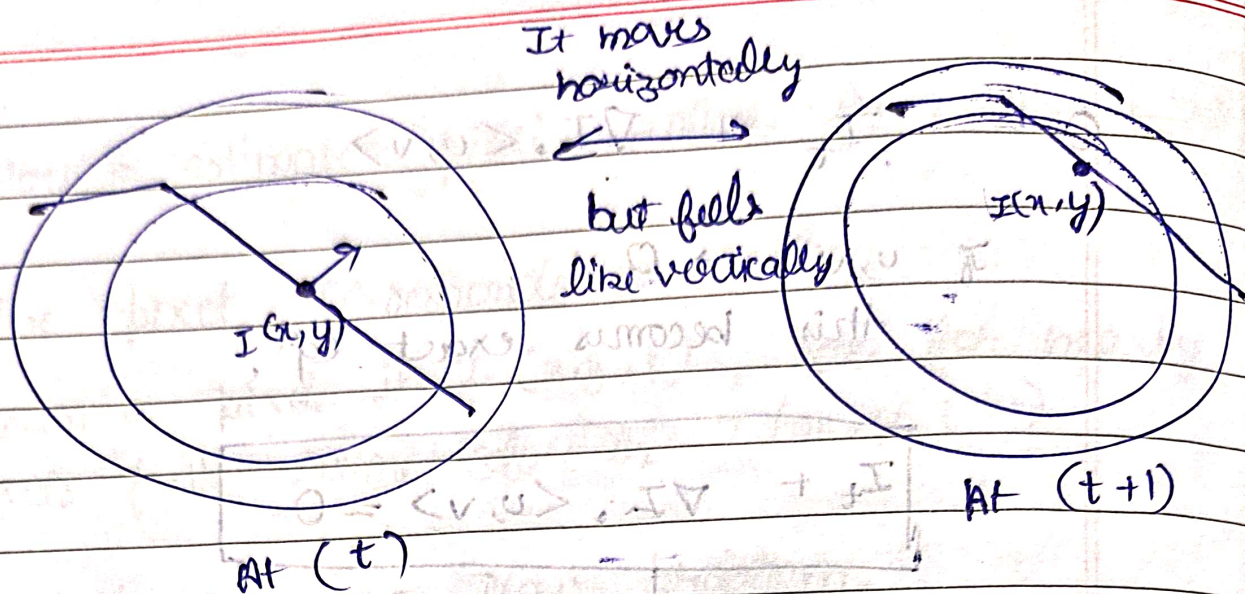
[So optical flow is "constrained" to be on a line]

Intuitively, this constraint means that

- The component of the flow in gradient direction is determined (called Normal Flow). I_x, I_y

- The component of the flow parallel to an edge is unknown

↓
called aperture problem



So, optical flow can't be computed from just brightness constancy constraint, we need additional constraints such as:

Global

Smooth optical flow (motion should be consistent in whole image.)

- Formulate error in optical flow constraint

$$e_c = \iint_{\text{img}} (I_x u + I_y v + I_t)^2 dx dy$$

Smoothness constraint

$$e_s = \iint_{\text{img}} (u_x^2 + u_y^2) + (v_x^2 + v_y^2) dx dy$$

Find (u, v) at each img. point that minimizes:

$$e = e_s + \lambda e_c$$

Local

eg.

1. The flow field is smooth locally

(Lucas-Kanade) ^{Method} Assume the pixel's neighbors have the same (u, v)

$$A d = -b$$

for a 5x5 window

25 eqⁿ, 2 unknown

do least squares & done

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

$$A d = -b$$

$$A^T A d = -A^T b$$

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

This is solvable when $A^T A$ is invertible

which means

λ_1 & λ_2 should not be too small

$\frac{\lambda_1}{\lambda_2}$ must not be too large
where $\lambda_1 > \lambda_2$

- Errors in Lucas-Kanade (B/c of assumptions)
- A point doesn't move like its neighbors we made
 - motion segmentation
- The motion is large (larger than a pixel) - Taylor doesn't hold.

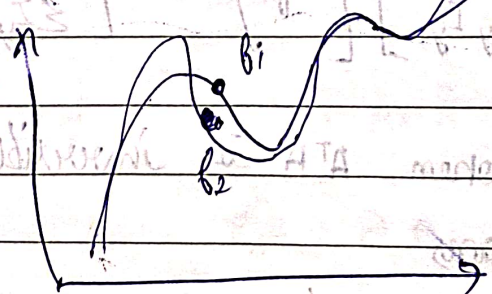
- Not-linear: Iterative Refinement
- Local minima: coarse to fine estimate

- Brightness constancy doesn't hold
 - exhaustive neighborhood search required.
- i) Iterative Refinement

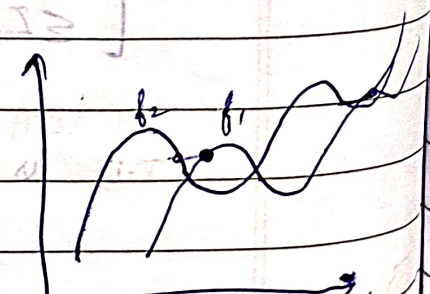
Algorithm

1. Estimate velocity at each pixel by solving Lucas-Kanade equations.
2. Warp I_t towards I_{t+1} using the estimated flow field.
3. Repeat until convergence.

II. Coarse to fine Refinement.



nearest match is correct (small motion)



nearest match is incorrect (aliasing)

Q. How to eliminate this problem of large motion?

Reduce the resolution until it moves only in range of 1 px motion.

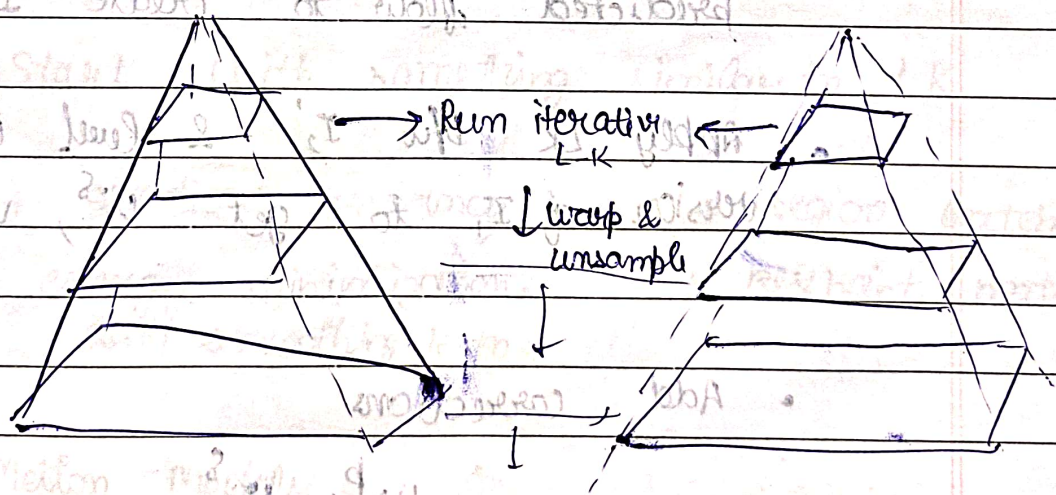


This is how it is done.

Hierarchical LK

i) Create a gaussian pyramid
At top, it is a blurry picture (coarse)

ii) Run iterative LK on the top level images.



iii) move downward, & repeat until you get the exact motion.

Before moving forward, let's take a detour to learn about aliasing & Image pyramids.

Formal Algorithm

1. Compute Iterative LK at level K
2. Init. $u_{K+1}, v_{K+1} = 0$ at size of level $K+1$
3. for each level $i = K$ to 0
 - upsample (Expand) u_{i+1}, v_{i+1} to create u_i^p, v_i^p flow field of new twice resolution as level i
 - Multiply u_i^p, v_i^p by 2 to get predicted flow
 - Work level i Gaussian version of I_2 w.r to predicted flow to create I_2'
 - Apply LK b/w I_2' & level i Gaussian version of I_1 to get u_i^s, v_i^s (correction in flow)
 - Add corrections

$$u_i = u_i^p + u_i^s$$

$$v_i = v_i^p + v_i^s$$

⇒ Sparse LK:

• The Lucas Kanade alg. described gives a dense field, (u, v) everywhere.

• But we said that we only want to solve LK where eigen-values are well behaved.

Sparse LK ⇒ hierarchical LK applied to good feature locations.

⇒ Today's State of the Art Methods:

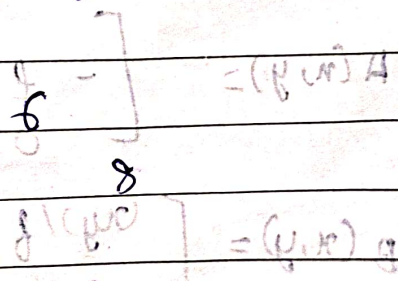
Start with something similar to LK

- + gradient consistency
- + region matching
- + energy minimization with smoothing term
- + keypoint matching

⇒ Motion Models

There are 4 type of motion in real world:

1. Translation
2. Un-known
3. Linear (similarity)
4. Affine
5. Homography/Perspective



Circular Motion

Now, in real world

$$V = \omega \times R + T$$

$$\begin{bmatrix} v_x \\ v_y \\ v_z \end{bmatrix} = \begin{bmatrix} 0 & -\omega_z & \omega_y \\ \omega_z & 0 & -\omega_x \\ -\omega_y & \omega_x & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} + \begin{bmatrix} v_{Tx} \\ v_{Ty} \\ v_{Tz} \end{bmatrix}$$

In image, velocity can be defined from

$$x = f \frac{X}{z}$$

$$y = f \frac{Y}{z}$$

$$\frac{dx}{dt} = v = v_x = f \frac{v_x}{z} - n \frac{v_z}{z}$$

$$\frac{dy}{dt} = v = v_y = f \frac{v_y}{z} - y \frac{v_z}{z}$$

$$\begin{bmatrix} v_x(n, y) \\ v_y(n, y) \end{bmatrix} = \frac{1}{z(n, y)} \begin{matrix} A(n, y) \cdot T + B(n, y) \cdot \Omega \\ \downarrow \qquad \qquad \downarrow \\ \text{Translation} \qquad \text{Rotation} \end{matrix}$$

$$A(n, y) = \begin{bmatrix} -f & 0 & n \\ 0 & -f & y \end{bmatrix}$$

$$B(n, y) = \begin{bmatrix} (ny)/f & -(f+y^2)/f & y \\ (f^2+y^2)/f & -(ny)/f & n \end{bmatrix}$$



Introduction to Tracking

Recognizing same object through frames in different timeline is called tracking.

Some great tracking ex.

i) surveillance camera

ii) Camera Tracking under occlusions

→ Tracking Challenges:

One idea to do tracking is optic flow b/w first & second, second & third & so so, but

i) It's hard to compute optical flow.

ii) There can be large displacements since could be moving rapidly

↓
Probably dynamics need to taken into account

iii) ~~Errors~~ would drift

iv) Occlusions

• Shi - Tomasi Feature Tracker

Algo

1. From frame t to frame $t+1$, track with LK and a pure translational model

- Robust for small disp.

2. Check consistency of tracks by affine registration to the first (or earlier) observed instance of the feature.

- Affine model is more accurate for larger displacement.

- Comparing to the first or early frame helps to minimize drift.

* This above method will fail for large disp.

• Tracking with dynamics

key idea: Given a model of expected motion, predict where objects will occur in next frame, even before seeing the image.

- Restrict search for the object
- Improved estimates since measurement noise is reduced by trajectory smoothness.

Assumptions

- Object do not disappear & reappear in diff. places in the scene
- Camera is not moving instantaneously to new viewpoint
- Gradual change in pose between camera and scene.

⇒ Tracking as inference:

Hidden state (x): True parameters we care about

Measurement (y): Noisy observation of underlying state.

At each time step t , state changes (from x_{t-1} to x_t) and we get a new observation y_t .

Our goal: Estimate distribution of state x_t given

→ all observations seen so far

→ Knowledge about dynamics of state transitions.



known

→ Steps of Tracking

i) Prediction: What is the next state of the object given past measurements?

$$P(x_t | y_0 = y_0, \dots, y_{t-1} = y_{t-1})$$

ii) correction: Compute an updated estimate of the state from prediction & measurements

$$P(X_t | Y_0 = y_0, \dots, \underbrace{y_{t-1}, y_t}_{\substack{\uparrow \\ \text{New} \\ \text{measurement}}})$$

Tracking \rightarrow The process of propagating this posterior distribution of state given measurements across time.

Some Assumptions

i) Only the immediate past matters

$$P(X_t | X_0, \dots, X_{t-1}) = P(X_t | X_{t-1})$$

\uparrow
Dynamics Model

ii) Measurements depend only on the current state

$$P(y_t | X_0, y_0, \dots, X_{t-1}, y_{t-1}, X_t) = P(y_t | X_t)$$

\uparrow
Observation Model

\rightarrow Tracking as Induction

Base Case: • Assume we have some initial prior that predicts state in the absence of any evidence: $f(X_0)$

• At the first frame, correct this, given value of $y_0 = y_0$

Prediction Step

$$\text{Given: } P(X_{t-1} | y_0, \dots, y_{t-1})$$

$$\text{Guess: } P(X_t | y_0, \dots, y_{t-1})$$

$$= \int P(X_t, X_{t-1} | y_0, \dots, y_{t-1}) dX_{t-1}$$

Law of Total Probability - Marginalization.

$$= \int P(X_t | X_{t-1}, y_0, \dots, y_{t-1}) P(X_{t-1} | y_0, \dots, y_{t-1}) dX_{t-1}$$

Using assumption

$$= \int \underbrace{P(X_t | X_{t-1})}_{\text{dynamics model}} \underbrace{P(X_{t-1} | y_0, \dots, y_{t-1})}_{\text{correct estimate from previous step}} dX_{t-1}$$

Correction Step

$$\text{Given: Predicted value } P(X_t | y_0, \dots, y_{t-1}) \text{ \& } y_t$$

$$\text{Compute } P(X_t | y_0, \dots, y_t)$$

predicted estimate

$$= P(y_t | X_t, y_0, \dots, y_{t-1}) P(X_t | y_0, \dots, y_{t-1})$$

$$P(y_t | y_0, \dots, y_{t-1})$$

(Using Bayes Rule)

$$P(y_t | X_t)$$

(Using assumption)

observed Model

→ Linear Dynamics model

Dynamics Model: state undergoes linear transformation plus Gaussian Noise

$$X_t \sim N(D_t X_{t-1}, \Sigma_{d_t})$$

↓ Normal Distribution ↓ Mean ↓ Gaussian Noise / Co-variance

→ Linear Measurement model

Observation Model: Measurement is linearly transformed state plus Gaussian Noise.

$$y_t \sim N(M_t X_t, \Sigma_{m_t})$$

⇒ Kalman Filter } gives a tradeoff b/w uncertainty of prediction & measurement.
 It is a method of tracking linear dynamic systems with Gaussian noise.

Kalman's Prediction

Mean at time $t =$ dynamics multiplier d
 *
 Corrected mean at time $t =$
 variance at $t =$ d^2 * corrected variance at $t =$
 + additional noise

μ_t^- , σ_t^- → Mean & covariance before measurement

μ_t^+ , σ_t^+ → vice versa

Kalman Correction (For 1-D)

$$\mu_t^+ = \frac{\mu_t^- \sigma_m^2 + m y (\sigma_t^-)^2}{\sigma_m^2 + m^2 (\sigma_t^-)^2}$$

$$(\sigma_t^+)^2 = \frac{\sigma_m^2 (\sigma_t^-)^2}{\sigma_m^2 + m^2 (\sigma_t^-)^2}$$

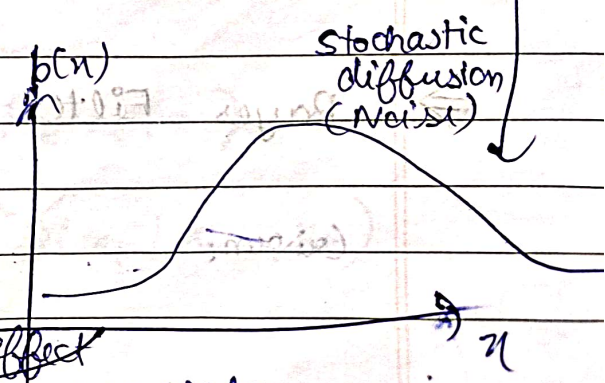
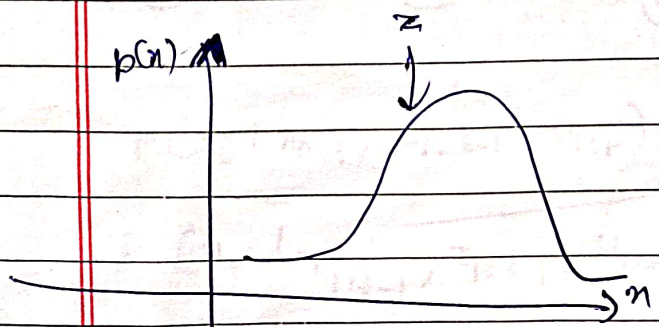
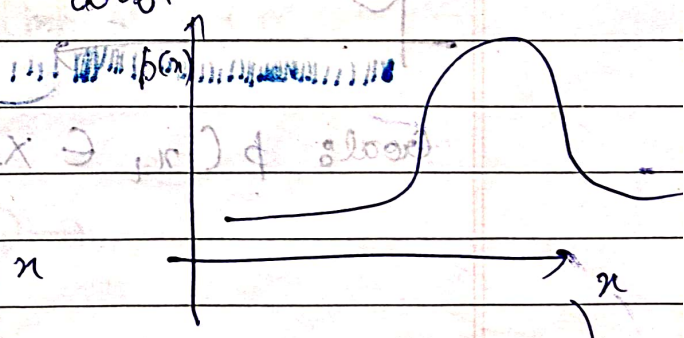
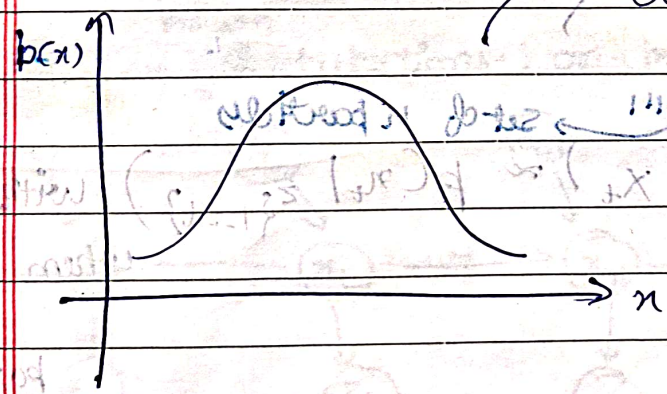
$$\mu_t^+ = \mu_t^- \left[+ K \left(\frac{y_1}{m} - \mu_t^- \right) \right]$$

Residual

* Also, variance always

Kalman Gain

Pictorial Representation



negative effect of measurement

Pros

- Simple updates, compact & efficient

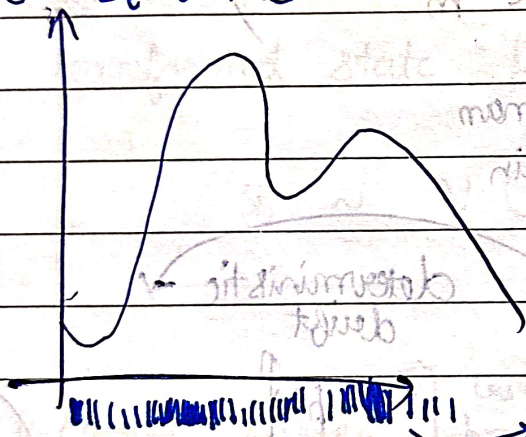
Cons

- Unimodal distribution, only single hypothesis
- Restricted class of motions defined by linear models + extensions called "extended KF Filter"

What about Non-Gaussian?

→ Particle Filters (YAY 😊)

Here $z_t \rightarrow$ measurements



Density is represented by both where the particles are and their weight.

→ set of n particles

Goal: $p(x_t \in X_t) \approx p(x_t | z_{1:t})$ with equality when $n \rightarrow \infty$

↓ particles

⇒ Bayes Filter

Concern:

1. Prior probability of the system state $p(x)$

2. Action (dynamical system) model:
 $p(x_t | u_{t-1}, x_{t-1})$

3. Sensor Model (likelihood) $p(z | x)$

4. Stream of observations z and action data u :

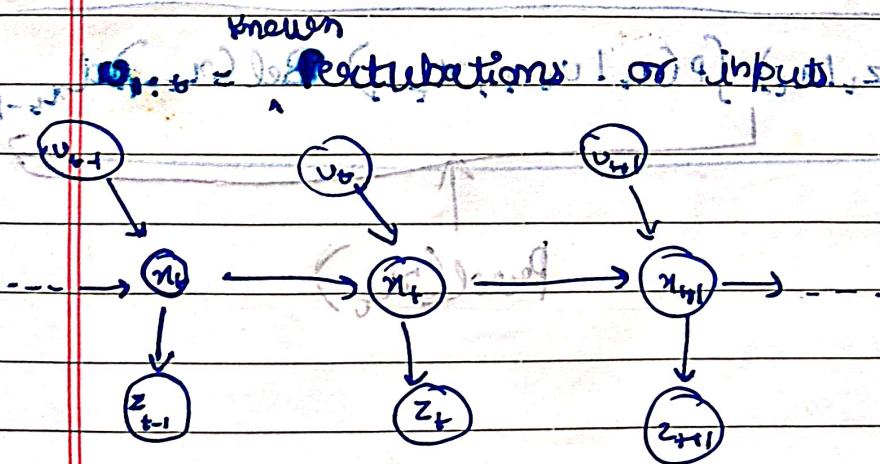
$$\text{data}_t = \{ u_1, z_1, \dots, u_{t-1}, z_{t-1} \}$$

Wanted:

1. Estimate of the state x at time t .

2. Posterior of the state is also called belief:

$$\text{Bel}(x_t) = P(x_t | u_1, z_1, \dots, u_{t-1}, z_{t-1})$$



$$P(z_t | x_{0:t}, z_{1:t-1}, u_{1:t}) = P(z_t | x_t)$$

$$P(x_t | x_{1:t-1}, z_{1:t-1}, u_{1:t}) = P(x_t | x_{t-1}, u_t)$$

→ Bayes Filter

$$\text{Bel}(x_t) = P(x_t | u_{1:t}, z_{1:t})$$

$$= \eta P(z_t | x_t, u_1, z_2, \dots, u_{t-1}) P(x_t | u_1, z_2, \dots, u_{t-1})$$

(Bayes)

$$= \eta \underset{\substack{\uparrow \\ \text{Sensor} \\ \text{Ind.}}}{P(z_t | x_t)} P(x_t | u_1, z_2, \dots, u_{t-1})$$

Sensor
Ind.

Total Prdb.

$$\text{of "Prior"} = \eta P(z_t | x_t) \int P(x_t | u_1, z_2, \dots, u_{t-1}, x_{t-1})$$

$$P(x_{t-1} | u_1, z_2, \dots, u_{t-1}) dx_{t-1}$$

Markov Ind.

$$= \eta P(z_t | x_t) \int P(x_t | u_{t-1}, x_{t-1})$$

$$\text{Bel}(x_{t-1}) dx_{t-1}$$

$$\text{Bel}(x_t) = \eta P(z_t | x_t) \int P(x_t | u_{t-1}, x_{t-1}) \text{Bel}(x_{t-1}) dx_{t-1}$$

