

# Lecture 6 (contd)

\* gaussian filter

give higher weight to pixels near the center

$$\frac{1}{273} * \begin{matrix} 1 & 4 & 7 & 4 & 1 \\ 4 & 16 & 26 & 16 & 4 \\ 7 & 26 & 41 & 26 & 7 \\ 4 & 16 & 26 & 16 & 4 \\ 1 & 4 & 7 & 4 & 1 \end{matrix}$$

2D gaussian  $g_g(x, y) = e^{-\left(\frac{x^2+y^2}{2a^2}\right)}$

Separable implementation

$$\begin{aligned} I_G &= I * g_g = \sum_i \sum_j I(i, j) e^{-\frac{i^2+j^2}{2a^2}} \\ &= \underbrace{\sum_i e^{-\frac{i^2}{2a^2}}}_{I * g_x} \underbrace{\sum_j I(i, j) e^{-\frac{j^2}{2a^2}}}_{g_y} \\ &= (I * g_x) * g_y = I * g_x * g_y \end{aligned}$$

Complexity of separable implementation

$M \times N$  image

$m \times m$  filter

One 2D pass:  $M \times N \times m^2$  operation

Two 1D pass:  $2 \times M \times N \times m$  operations

$$2(m) \cdot MN < MN(m)^2$$

repeated application of gaussians:

$$I * G_{\sigma_1} * G_{\sigma_2} = I * G_{\sqrt{\sigma_1^2 + \sigma_2^2}}$$

$$2 \times (M \times N \times m^2) < (M \times N \times (2m)^2)$$

? Selecting gaussian variance ( $\sigma_a$ )

high the sigma  $\rightarrow$  bigger the filter

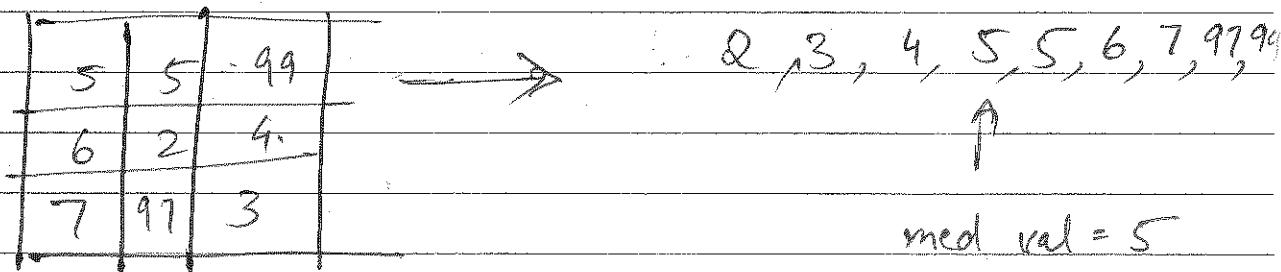
rule is to use filter  $m=5a \Rightarrow a \leq \frac{m}{5}$

$$\frac{1}{2 \cdot 5} \boxed{0.14 | 0.01 | 0.61 | 0.14} \quad a = \frac{5}{6}$$

$$\frac{1}{16} \boxed{1 | 4 | 6 | 4 | 1} \quad a = 1$$

\* Median filter :

$$I_{\text{med}}(i, j) = \text{median} \{ I(x, y) | (x, y) \in \text{with}(i, j) \}$$



good for salt & pepper noise.

\* Multiple scale analysis:

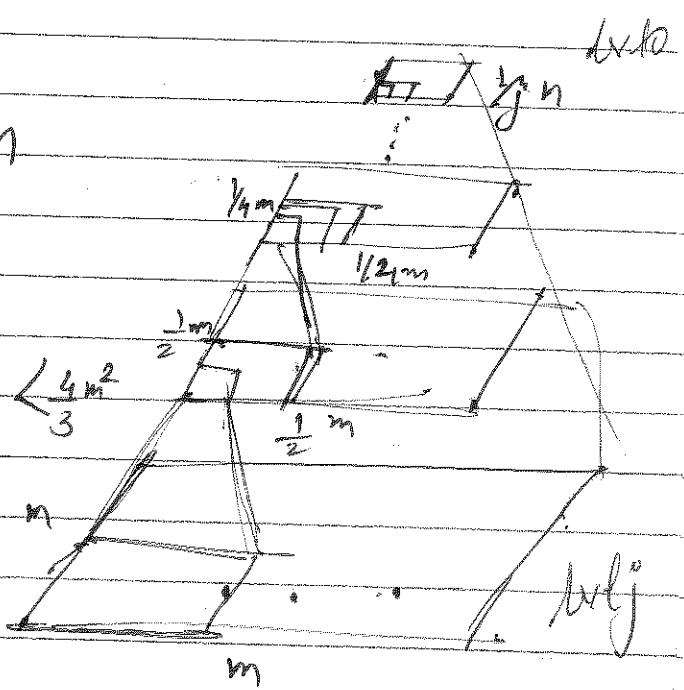
analyse image with different filter size

Image pyramid:

$$m = 2^J \Rightarrow J = \log_2 m$$

total # pixels =

$$m^2 + \frac{1}{4} m^2 + \frac{1}{16} m^2 + \dots < 4m^2$$



Gaussian pyramid:

$$I_j \rightarrow [G] \rightarrow [\downarrow 2] \rightarrow I_{j-1}$$

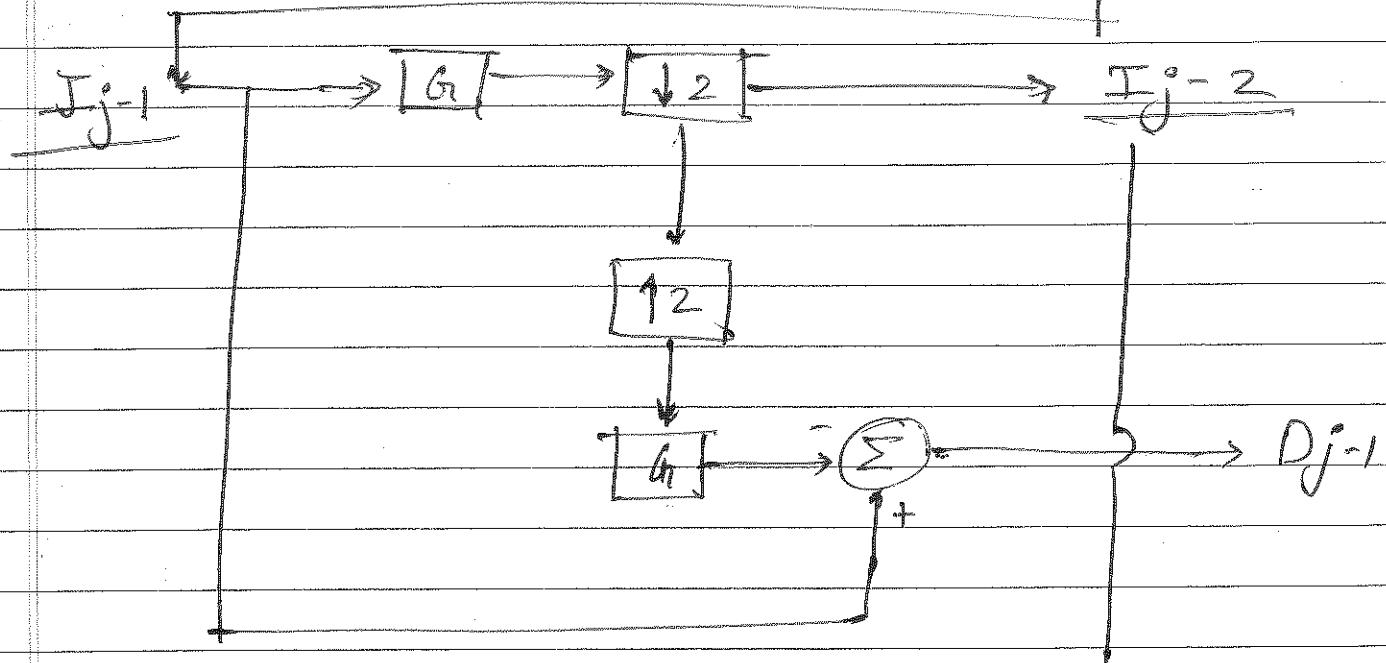
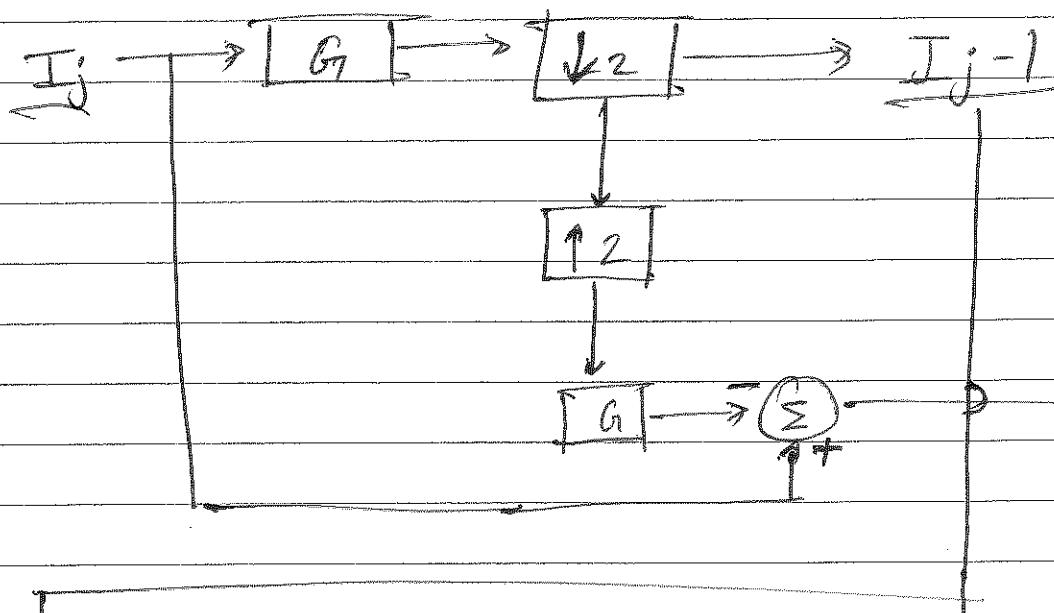
$$I_{j-1} \rightarrow [G] \rightarrow [\downarrow 2] \rightarrow I_{j-2}$$

$$I_{j-2} \rightarrow [G] \rightarrow [\downarrow 2] \rightarrow I_{j-3}$$

- Helps to tackle Aliasing images
- useful for analysis

## ~~Skew~~ Laplacian pyramid:

- more use for image compression
- variance is smaller



CLASSMATE

Date \_\_\_\_\_  
Page \_\_\_\_\_

# Lecture 7

## ★ Feature Detection

anything like eyes, hair, etc.

Two approach:

1. Feature engineering
2. Deep learning

### 1. Feature engineering

used when problem is known or available

\* feature types

- a) global: eg color distribution (histogram)
- b) local: eg edges or corners

\* goal:

- characterize images  
and/or
- local neighborhoods

#### 1.1. Edge detection

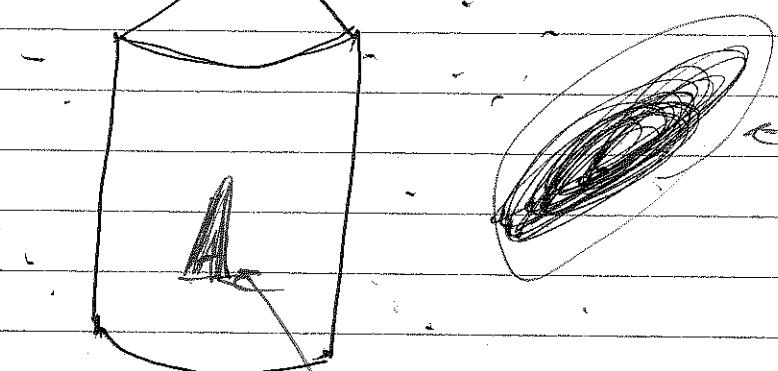
requirements:

- correspond to scene elements
- invariant [illumination, pose, viewpoint, scale]
- reliable detection

edge  $\Rightarrow$  location with change in image

depth discontinuity

normal discontinuity



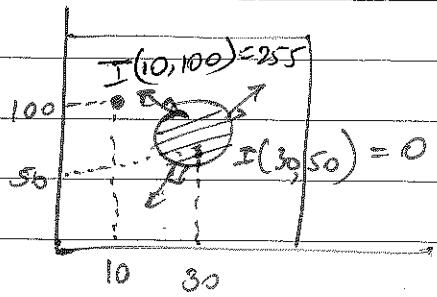
reflection (or color) discontinuity

### \* edge detection steps:

- 1) Smooth image to reduce noise (without affecting edges.)
- 2) Enhance edges
- 3) Detect edges
- 4) Localize edges.

## Image gradient

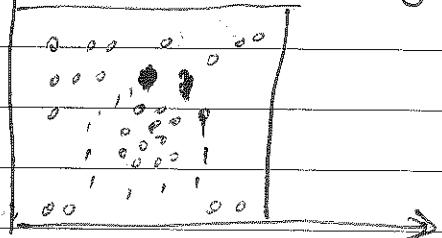
Image :  $I(x,y)$



$$VI(x,y) = \begin{bmatrix} \frac{\partial I}{\partial x} \\ \frac{\partial I}{\partial y} \end{bmatrix}$$

↑ (i) zero when no charge

$$= \begin{vmatrix} I_x \\ I_y \end{vmatrix}$$



Note: image gradient vectors are  $\perp$  to the edges.

$$\frac{\partial I}{\partial x} = \lim_{h \rightarrow 0} \frac{I(x+h, y) - I(x, y)}{h}$$

$$\frac{\partial I}{\partial y} = \lim_{h \rightarrow 0} \frac{I(x, y+h) - I(x, y)}{h}$$

$\rightarrow$

$$\frac{\partial}{\partial x} I(x,y) = I(x+1,y) - I(x,y)$$

edge  $\Rightarrow$  edge element

$$\text{magnitude} : |\nabla I| = \sqrt{x^2 + y^2}$$

$$\text{angle} : \theta = \tan^{-1}\left(\frac{I_y}{I_x}\right)$$

\* Edge detection :

$$E(i,j) = \begin{cases} 1 & \text{if } |\nabla I(i,j)| > \tau \\ 0 & \text{otherwise} \end{cases}$$

Threshold

Edgemap

• SO,

$$\frac{\partial I(x,y)}{\partial x} = I(x+1, y) - I(x, y)$$

$$\Delta x = \begin{bmatrix} (-1) & 0 & 1 \\ 0 & (-1) & 0 \\ 1 & 0 & (-1) \end{bmatrix}$$

$$y \rightarrow \begin{bmatrix} +1 & +1 \end{bmatrix}$$

x → z + 1

Similarly,

detects vertical edges

$$\Delta y = \begin{bmatrix} 1 & 0 & 0 \\ 0 & (-1) & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

## \*Central difference

$$\frac{\partial}{\partial I} I(x, y) = \lim_{h \rightarrow 0} \frac{I(x+h, y) - I(x-h, y)}{2h}$$

here  $h=1$

$$\begin{bmatrix} -1 & 0 & 1 \end{bmatrix}_{x-1 \quad x \quad x+1} = I(x+1, y) - I(x-1, y)$$

$$\Delta_x = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} \xrightarrow{\text{derivative}} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix} \xrightarrow{\text{smooth}} \Delta_y = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix} \xrightarrow{\text{desi}}$$

Note : Explanation,

$$\frac{\partial}{\partial I} I(x, y) = (+) I(x+1, y) + 0 - I(x-1, y)$$

$$(+ \underset{\downarrow}{x+1} + 0 - \underset{\downarrow}{x-1})$$

$$(-1 \quad 0 \quad 1)$$

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

\* Implementation using -

\* ~~Sobel filter~~ : (Open CV)

$$\Delta x = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} * \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

noise reduction      edge detection

$$\Delta y = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

Example 6

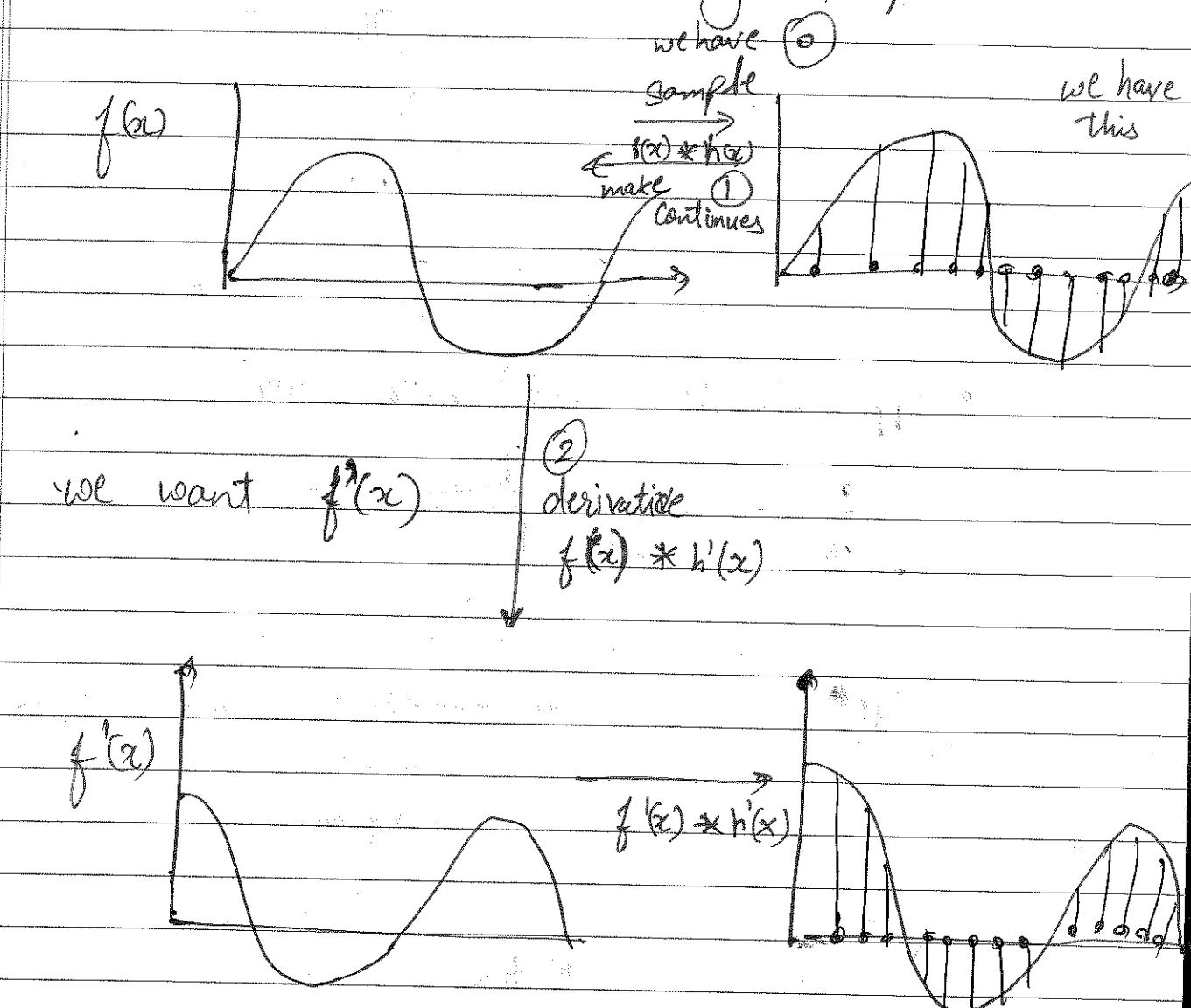
$$\begin{array}{|c|c|c|c|c|} \hline & \begin{bmatrix} -1 & 1 \\ -2 & 2 \end{bmatrix} & 255 & 255 & 255 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} = 1 \times 255 + 2 \times 255 + 1 \times 255 = 1020$$

$$\Delta x = \begin{bmatrix} 0 & 0 & 1020 & 1020 & 0 & 0 \\ 0 & 0 & 1020 & 1020 & 0 & 0 \\ 0 & 0 & 1020 & 1020 & 0 & 0 \\ 0 & 0 & 1020 & 1020 & 0 & 0 \\ 0 & 0 & 1020 & 1020 & 0 & 0 \end{bmatrix}$$

$$\Delta y = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ -1 & -1 & -1 & -1 & -1 \end{bmatrix}$$

## \* Implementation Using -

- More accurate derivative (using 1D) explained :

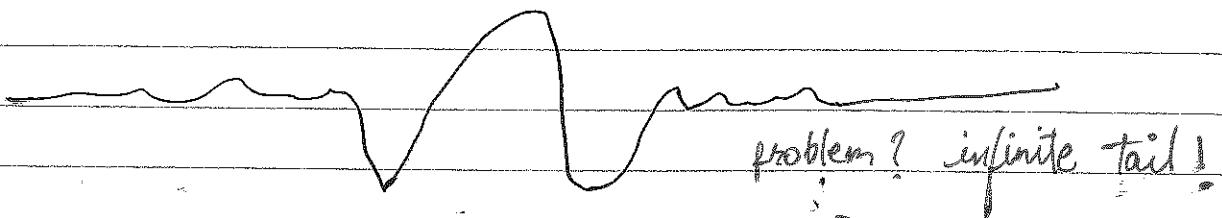


So we have to define  $h(x)$ , then convolve the continuous func with derivative of  $h(x)$  [ $h'(x)$ ].

- Perfect interpolator

Under certain conditions perfect reconstruction.

$$h(x) = \text{sinc}(x) = \frac{\sin(\pi x / T)}{(\pi x / T)}$$



- approximate  $h(x)$  with  $G(x)$

$$f_s \geq 2 f_m$$

sampling frequency  
(pixel/mm)

variation

less variation then less  $f_m$

$\uparrow f_m = \uparrow$  variation

$\downarrow f_m = \downarrow$  variation

$$I_x = I * G_x$$

$$I_y = I * G_y$$

$$f[x] * g[x]$$

using separable property of gaussian

$$I_x = I * G[x] * G[y]$$

$$I_y = I * G[y] * G[x]$$

$$I_x = \boxed{\text{image}} * \boxed{\text{derivative}} \boxed{\text{gaussian}} * \boxed{\text{smooth}}$$

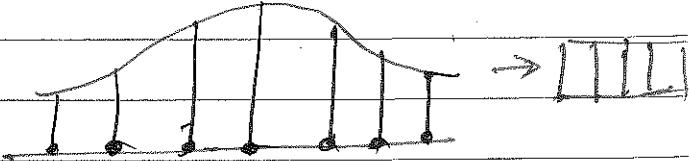
$$\cancel{I_x} = I * G_x[x] * G_y[y]$$

$$I_y = I * G_y[y] * G_x[x]$$

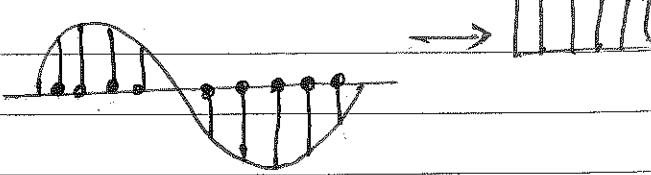
$$I_y = \boxed{\text{image}} * \boxed{\frac{d}{dx}} * \boxed{\text{smooth}} \boxed{\text{gaussian}}$$

So 1D gaussian,

$$G_x = e^{-\frac{x^2}{2\sigma^2}}$$

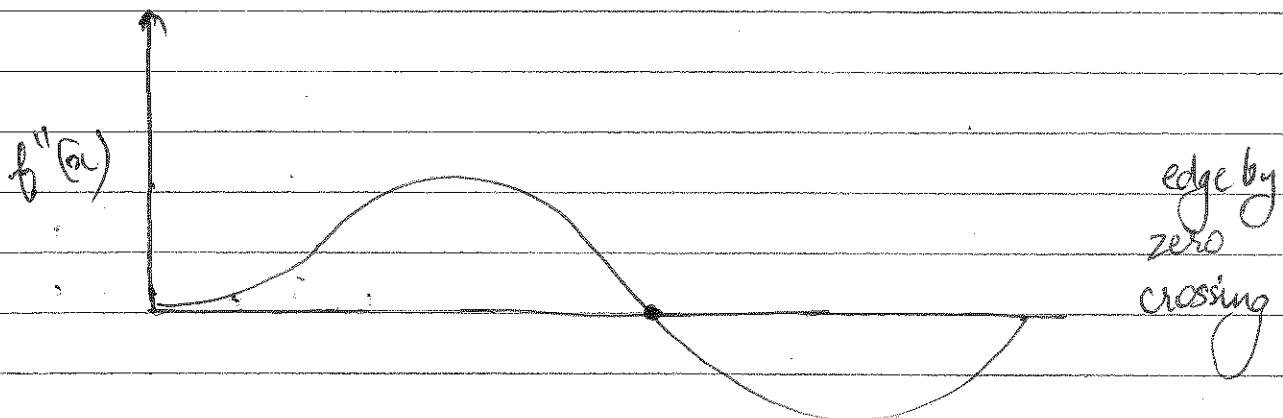
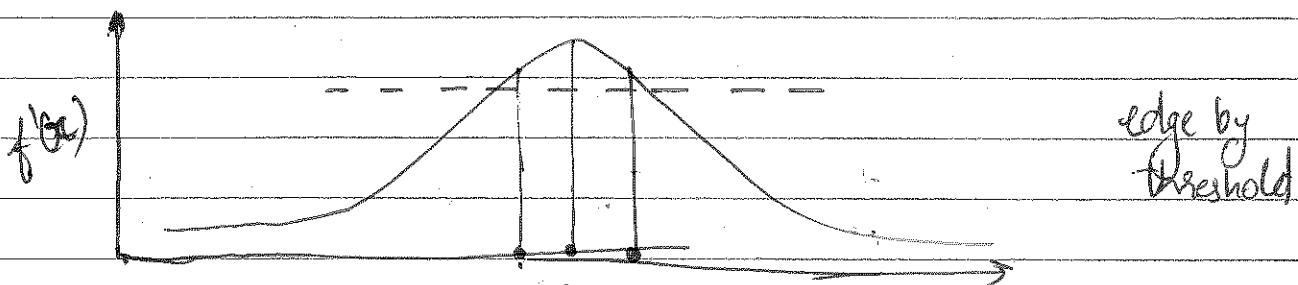
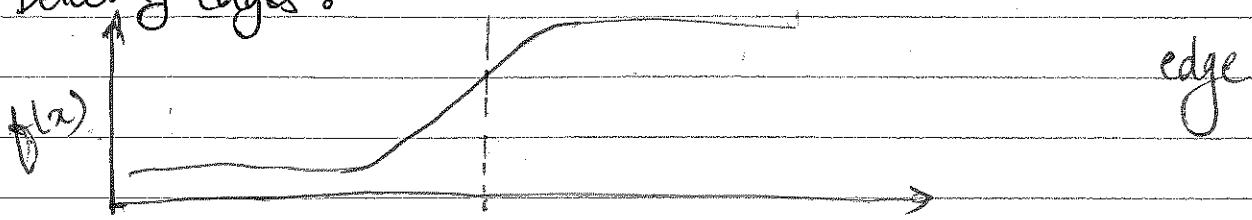


$$G'_x = -\frac{2x}{2\sigma^2} e^{-\frac{x^2}{2\sigma^2}}$$



# Lecture - 8 (feature detection contd.)

Detecting edges :



\* In 2D, second derivative as scalar quantity:

$$\text{Laplacian : } \nabla^2 I = \frac{\partial^2 I}{\partial x^2} + \frac{\partial^2 I}{\partial y^2} = I_{xx} + I_{yy}$$

~~$I_x = I(x+1) - I(x)$~~

$$I_{xx} = (I(x+1) - I(x)) - (I(x) - I(x-1))$$

$$= I(x+1) - 2I(x) + I(x-1)$$

$$\Rightarrow \boxed{1 \quad -2 \quad 1}$$

$I_{yy} \Rightarrow$ 

1
-2
1

 $I_{xx} + I_{yy} \Rightarrow$ 

0	1	0
1	-4	1
0	1	0

 $\rightarrow$  Laplacian filter

Apply this and look for crossing edges.

Smooth with gaussian before applying laplacian (bcz 2nd degree derivative is more sensitive to noise)

$$H = \nabla^2(I * G) = \nabla^2 G * I$$

$\downarrow$

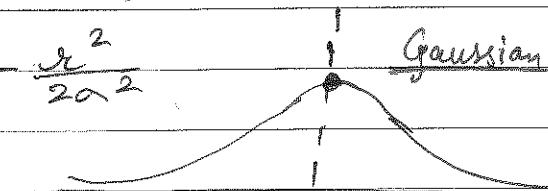
(filter)  $\text{LOG} \rightarrow$ 

Lap of Gaussian

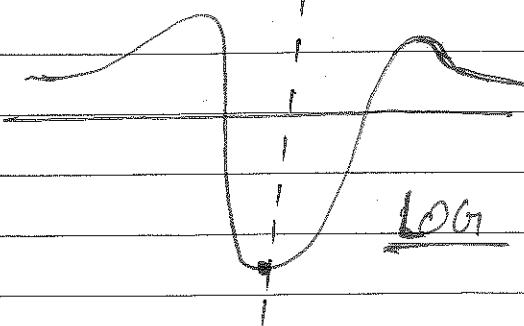
we know,

$$g = e^{-r^2/2\sigma^2} \quad (r^2 = x^2 + y^2)$$

$$\nabla^2 g = \frac{r^2 - 2\sigma^2}{r^4} e^{-\frac{r^2}{2\sigma^2}}$$



Gaussian



LOG

\* Edge detection using LOG :

1) Compute  $\text{LOG}^\circ : H = (\nabla^2 G) * I$

2) Threshold :  $E(i,j) = \begin{cases} 0 & \text{if } H(i,j) < 0 \text{ (output a bin)} \\ 1 & \text{if } H(i,j) \geq 0 \text{ (image)} \end{cases}$

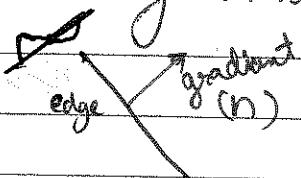
3) Mark the (zero crossing) edges at threshold  
(ie  $0 \rightarrow 1$  or  $1 \rightarrow 0$ )

Log  $\rightarrow$  Image =  $\frac{1 - \min}{\max - \min} \times 255$

\* Canny edge detection :

Detect edges at zero crossings of second order  
diff derivatives taken along the gradient. (direction derivative)

$n$  = gradient



if  $|n| > T$  detect edges at zero crossing  
of  $\frac{\partial^2}{\partial n^2} (I * G)$

\* Comments :

1) attempts detect only if gradient magnitude is large enough ( $|n| > T$ )

2) smooth along edges to preserve edges.

3) Alternative to zero crossing of  $\frac{\partial^2}{\partial n^2} (I * G)$  is

maximum at  $\frac{\partial}{\partial n} (I * g)$ .

- so let's see ~~how it does~~ what we have to do,

gradient:  $n = \nabla (I * g)$

directional derivative:  $\frac{\partial (I * g)}{\partial n}$  =  $n \cdot \nabla (I * g)$

↑ direction      ↑ gradient

$$n \cdot \nabla I * g = \nabla (I * g) \cdot \nabla (I * g)$$

↑

direction of gradient

$$= |\nabla (I * g)|^2$$

↑

mag of gradient

So we find max of this and that our edge.

### Summary:

if  $|\nabla (I * g)| > T$  set edge at max  
if  $|\nabla (I * g)|$  along the direction of  $\nabla (I * g)$

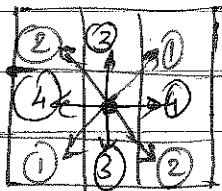
### \* Additional Components:

#### 1) Non-maximum suppression

Need: local maximum of gradient mag in direction of gradient. (local max)

$$\nabla(I * h) = (I_x, I_y)$$

$$\theta = \tan^{-1} \left( \frac{I_y}{I_x} \right)$$



$$\theta^* = \text{round} \left( \frac{\theta}{45} \right) * 45$$

0, 45, 90, 135...

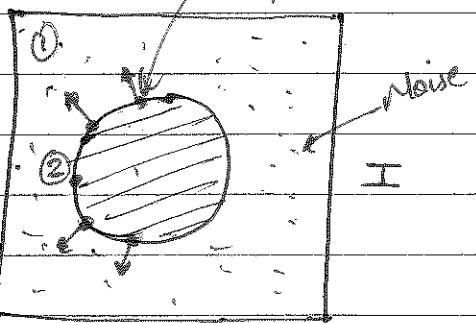
$$E(i,j) = \begin{cases} 1 & \text{if } \nabla(I * h) \text{ is a local max along} \\ & \text{the edge} \\ 0 & \text{otherwise} \end{cases}$$

## 2) hysteresis thresholding (tracking)

use  $Z_H$  to start tracking and  $Z_L$  to continue ( $Z_H > Z_L$ )

edge points

(a) Initialize array of visited pixels  $V(i,j) = 0$



(b) Scan Image T-B, L-R:  
if  $|V(i,j)| < |V| > Z_H$ ,  
start tracking on edge.

(c) Search for additional neighbors in direction orthogonal to  $\nabla V$  such that  $|V| > Z_L$

\* Edge detection performance evaluation:

1) Theoretical

2) Experimental (true positive, false positive)

3) Application outcome

\* Summary of (1.1.) Edge detection

Edge detection requires finding discrete image derivatives:

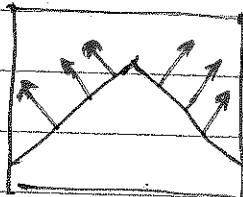
- forward difference
  - central difference
  - Sobel
  - Gaussian derivatives
  - zero crossing of LOG (2nd order)
  - Canny (directional derivative)
- } first order derivatives

## 1. Feature Detection:

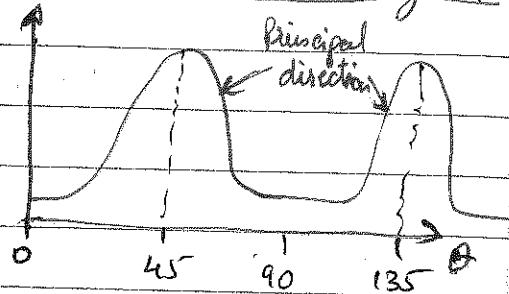
## 1.2 CORNERS

corner = more than one direction in orientation histogram.

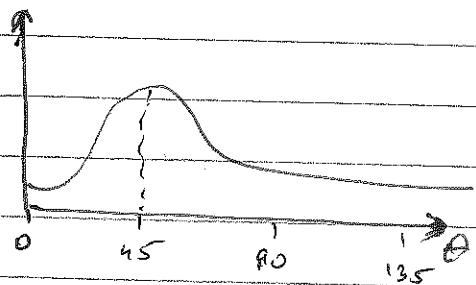
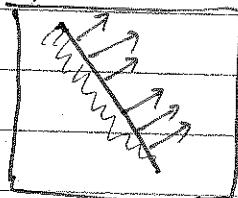
corner



orientation histogram



edge



A corner is better defined than edge, it has a location. (better localization)

## \* Algorithm outline:

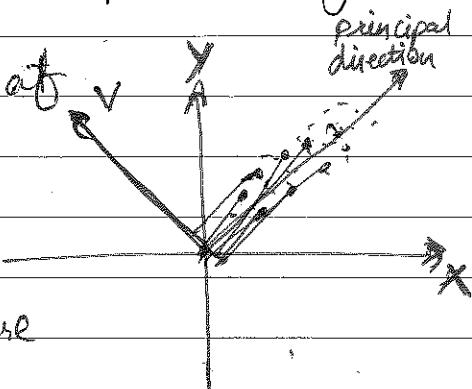
① find correlation matrix of gradients in local ~~matrix~~ window.

② find eigen values at correlation matrix  $(2 \times 2)$

③ detect corner in window if eigenvalues are sufficiently large.

\* Algorithm is based on principle component Analysis (PCA)

Goal: find direction  $v$  s.t. projection of  $v$  onto  $\{g_i\}$  is minimized.



① Each  $\{g_i\}$  has  $\begin{pmatrix} x_i \\ y_i \end{pmatrix}$  and are a point in coordinate system.

② we try to find a  $v$  such that the points  $\{g_i\}$  once projected on  $v$  minimize the spread.

So, we use  $E(v)$  where  $E(v)$  tells us how good the direction is.

$$E(v) = \sum (g_i \cdot v)^2 = \sum (g_i^T v)(g_i^T v)$$

$(g_i \cdot v)$  gives projections of  $\{g_i\}$  on  $v$ .

$$= \sum_i (v^T g_i)(g_i^T v) = \sum_i v^T g_i g_i^T v$$

$$E(v) = v^T (\sum_{g_i} g_i g_i^T) v = v^T C v$$

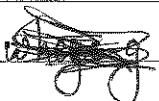
$1 \times 2 \quad 2 \times 2 \quad 2 \times 1$

$1 \times 1$

$\{g_i\}$  is a  $2 \times 1$  vector & so  $C$  is  $2 \times 2$  matrix.

$C$  is called correlation matrix.

$$C = \sum_i g_i g_i^T = \begin{bmatrix} \sum x_i^2 & \sum x_i y_i \\ \sum x_i y_i & \sum y_i^2 \end{bmatrix}$$



Note: In python we stack  $\{g_i\}$  on top of each other to form a data matrix  $D$  & find correlation matrix.

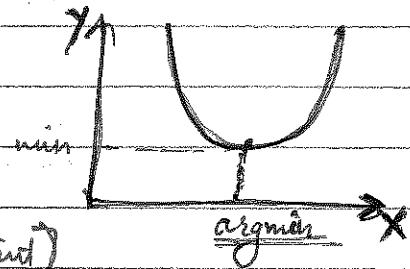
$$D = \begin{bmatrix} g_1^T \\ g_2^T \\ \vdots \\ \vdots \\ g_m^T \end{bmatrix}$$

$$\text{Q} \& C = D^T D$$

Note

so,

$$\left\{ \begin{array}{l} E(v) = v^T C v \\ v^* = \arg \min_v E(v) \end{array} \right. \quad C \in \mathbb{R}^{m \times m}$$



$\nabla E(v) = 0$  — [is the minimum point]

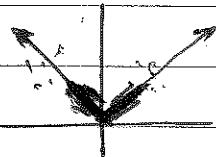
i.e.  ~~$Cv = 0$~~   $\Rightarrow$

solution is eigen vector belonging to smallest eigenvalues.

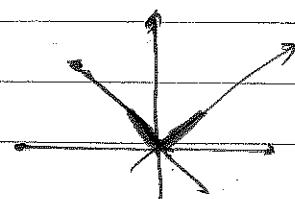
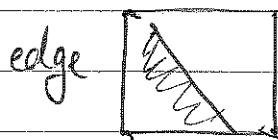
$$\boxed{C}_{2 \times 2} \boxed{v} = \boxed{0}_2$$

eigenvalue = variance in corresponding principal direction

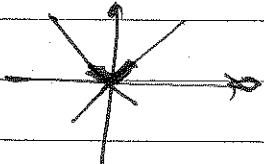
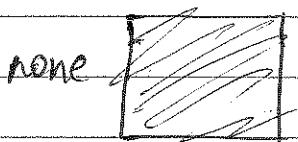
# Lecture 9 - W5L1



$\lambda_1$  large  
 $\lambda_2$  large



$\lambda_1$  small  
 $\lambda_2$  large



$\lambda_1$  small  
 $\lambda_2$  small

Corner if  $\lambda_1 \cdot \lambda_2 > \tau$



Corner detection Summary

(1) We scan the image T-B, L-R, at each pixel we select a local  $3 \times 3$  neighborhood.

(2) Build correlation matrix  $C = \sum g_i g_i^T$

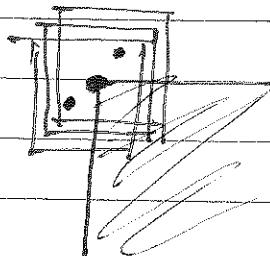
(3) Compute eigenvalues of  $C$

(4) Detect corner if  $\lambda_1 \cdot \lambda_2 > \tau$

### \* Non - maximum suppression

If not done , each corner will lead to detecting multiple points.

So to get proper point we will need to :



(1) compute  $\lambda_1, \lambda_2$  for all windows.

(2) ~~Select~~ Select  $\lambda_1, \lambda_2 > T$  windows and sort in descending order.

(3) Select the top of the list as corner , and delete all other corners in its neighborhood.

(4) Stop once detecting X% of the points as corners.

This analysis needs to be done at multiple scales.

\* Harris corner detection :

(1) Compute correlation matrix  $C$  for window.

(2) Compute cornerness measure :

$$G(c) = \frac{\det(C)}{\lambda_1 \cdot \lambda_2} - \frac{K \operatorname{tr}^2(C)}{K(\lambda_1 + \lambda_2)^2}$$

$K$  is selected  
By user

(3) detect corners where  $G(c)$  is high

! trace is ~~sum~~ sum of diagonal elements

$K \in [0, 0.5]$  is a user parameter

if  $K = 0$   $G(c)$  detects corners

if  $K = 0.5$   $G(c)$  detects edges

So,

$$G(c) = \det(C) - K \operatorname{tr}^2(C)$$

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

$$= (1-2K) \lambda_1 \lambda_2 - K(\lambda_1^2 + \lambda_2^2)$$

= 0 if  $K = 0.5$

corner detection

= 0 if  $K = 0$

edge detection

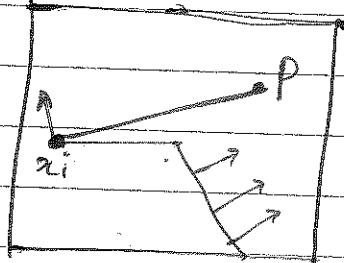
Alternative basis:

$$C(c) = \det(c) + k \left( \frac{\text{tr}(c)}{2} \right)^2$$

\* Corner localization:

Given that there is a corner in a window find its location.

To determine if  $p$  is the corner connect each point  $x_i$  to  $p$  and project the gradient at  $x_i$  onto  $(x_i - p)$



The "best"  $p$  will minimize the sum of all projection

$$E(p) = \sum_i (\nabla I(x_i) \cdot (x_i - p))^2 \quad [0 \text{ if corner}]$$

$$= \sum_i (x_i - p)^T \nabla I(x_i) \nabla I(x_i)^T (x_i - p)$$

$$= \sum_i (x_i - p)^T (\nabla I(x_i) \nabla I(x_i)^T) (x_i - p)$$

$$\begin{matrix} & (2 \times 1) & (1 \times 2) \\ \underbrace{1 \times 2} & & \underbrace{2 \times 2} & \underbrace{1 \times 1} \end{matrix}$$

$$\left\{ \begin{array}{l} E(p) = \sum_i (x_i - p)^T (\nabla I(x_i) \nabla I(x_i)^T) (x_i - p) \\ p^* = \underset{p}{\operatorname{arg\min}} E(p) \\ \Rightarrow \nabla E(p) = 0 \end{array} \right.$$

$$\cancel{-2 \sum_i (\nabla I(x_i) \nabla I(x_i)^T) (x_i - p)} = 0$$

$$\sum_i \nabla I(x_i) \nabla I(x_i)^T p = \underbrace{\sum_i \nabla I(x_i) \nabla I(x_i)^T}_{C \quad 2 \times 2} \underbrace{x_i}_{2 \times 1} \quad \frac{1}{2x_1} \quad \frac{2 \times 2}{2x_1}$$

To solve it, isolate  $p$  on one side & move everything on other side.

$$CP^* = \cancel{\sum_i \nabla I(x_i) \nabla I(x_i)^T} x_i$$

multiply by  $C^{-1}$

~~$\nabla I(x_i)$~~

$$P^* = \left( \sum_i \nabla I(x_i) \nabla I(x_i)^T \right)^{-1} \sum_i \nabla I(x_i) \nabla I(x_i)^T x_i$$

$$P^* = C^{-1} \sum_i \nabla I(x_i) \nabla I(x_i)^T x_i$$

location  
of corner

correlation  
matrix.

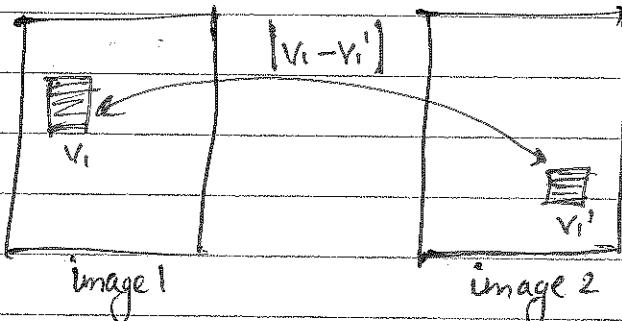
## \* Feature point characterization

Given a corner we would like to compare it to corners in another image and to use it to characterize the image.

⇒ feature point characterization (find feature vector)

### \* Applications:

- matching
- tracking
- correspondence
- recognition



### \* Example methods:

- HOG
- SIFT
- SURF

- Shape context
- spin image

### \* Desired Properties

- translation invariance → local window
- rotation invariance → histograms
- scale invariance → pyramid
- illumination invariance → gradients of gradients

⇒ orientation histograms at different scales

+ rotation invariance

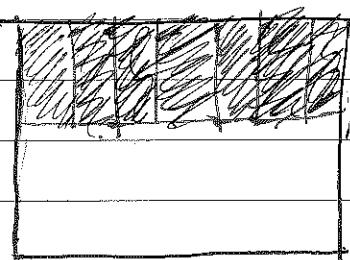
Rotate coord system (rotate hist) so that dominant direction at start.

+ Scale invariance

Need representation at multiple scales.

### \* HOG method: (OpenCV)

(1) split each patch into cells  
(possibly overlapping)



(2) create orientation histogram  
in each cell (using edge or gradient directions,  
possibly weighted by distance from center of  
gradient magnitude)

(3) Concatenate orientation histograms

### \* SIFT (Scale Invariant Feature Transform):

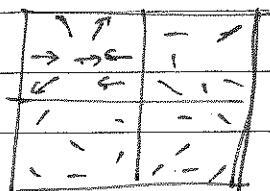
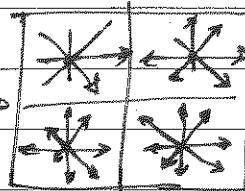


Image gradients



Keypoint descriptor

- ⇒ - use weighted sum to create orientation histograms in cells, then concatenate
- Align histogram based on dominant direction (rotation invariance).

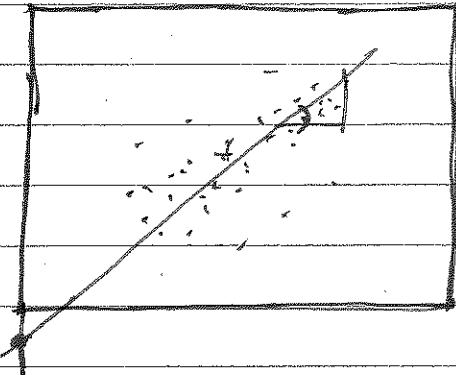
# Lecture 10 - WSL2

## MODEL FITTING

- Line detection Fitting

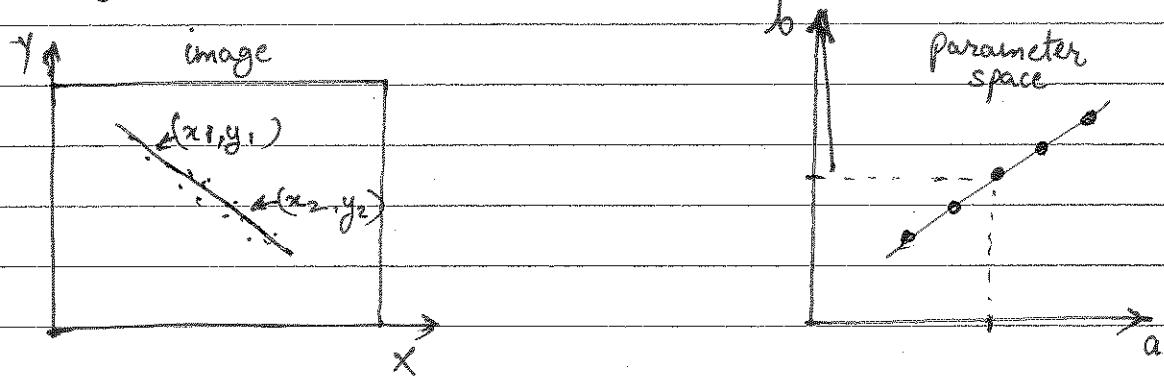
Tasks:

- 1) Grouping
- 2) fitting



Two simultaneous tasks. Solving one makes the other easier.

\* Hough transform:



$$y = ax + b$$

\* multiple points on a line in the image correspond to a single point in parameter space (the parameters of the line).

\* Detect lines by casting votes in parameter space.

$$\text{so, } y = ax + b \rightarrow b = y - ax$$

given  $(x, y)$  scan a and compute b. cast votes at locations  $(a, b)$  - a line.

- Parameters with highest votes wins in parameter space. Intersections in parameter space gives parameters common to multiple points in the image space.

- When doing this on image votes need to be stored in a matrix. What so should be the intercepts of parameter space ie slope & y-intercept.

slope can't go to  $\infty$  so it can be problematic.  
So that cannot be a good model.

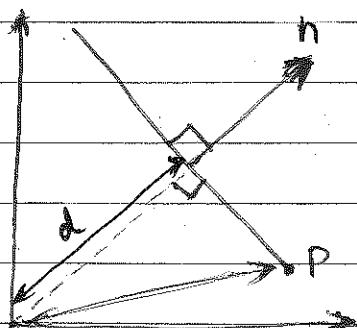
→ Better line representation:

\* Problem with  $y = ax + b$  model:

- range of  $a$ ?

- How to represent vertical lines?

Solution:



$n = (nx, ny)$  normal  
(orientation)

$d$  = distance from origin

for point  $p = (x, y)$  to be on the line:

$$(n_x, n_y) \cdot (x, y) = d \quad \leftarrow \text{line equation using } n, d$$

$$\therefore n \cdot p = d$$

$$\text{where } \| (n_x, n_y) \| = 1$$

\* This is called as **Implicit line equation**:

$$n = (n_x, n_y) \quad \text{where } \begin{cases} n_x = \cos \theta & \text{angle of} \\ n_y = \sin \theta & \text{normal} \end{cases}$$

equation is,

$$(n_x, n_y) \cdot (x, y) = d$$

$$\Rightarrow (\cos \theta)x + (\sin \theta)y - d = 0$$

$$(Ax + By + c = 0)$$

$\underbrace{\qquad\qquad\qquad}_{\text{normal}}$  negative dist from origin

line parameters:  $\theta, d$

\* Hough with explicit line equation:

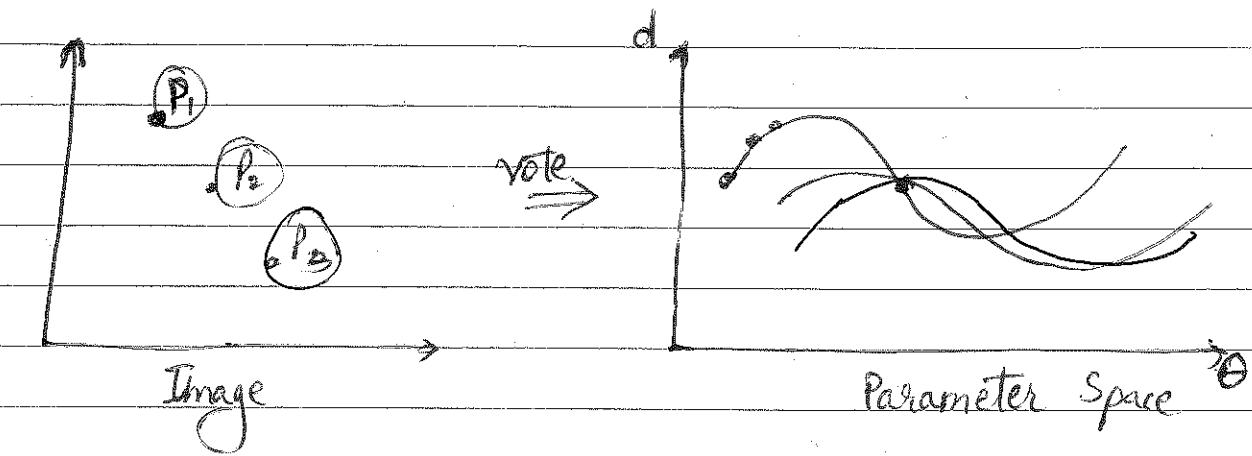
Implicit

$\rightarrow$  ~~Explicit~~ line equation:

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

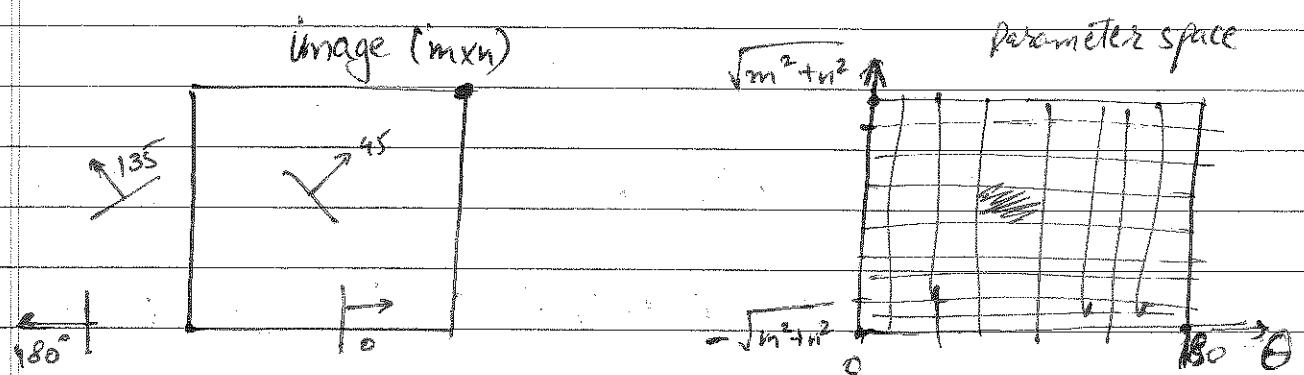
$\rightarrow$  Vote by point  $(x, y)$ :

for each  $\theta_i$  vote  $d_i = x \cos \theta_i + y \sin \theta_i$



## \* Practical issues

### 1) Parameter space size



2) Bin size : larger bins are more efficient but provide less localization. Small bins are more accurate.

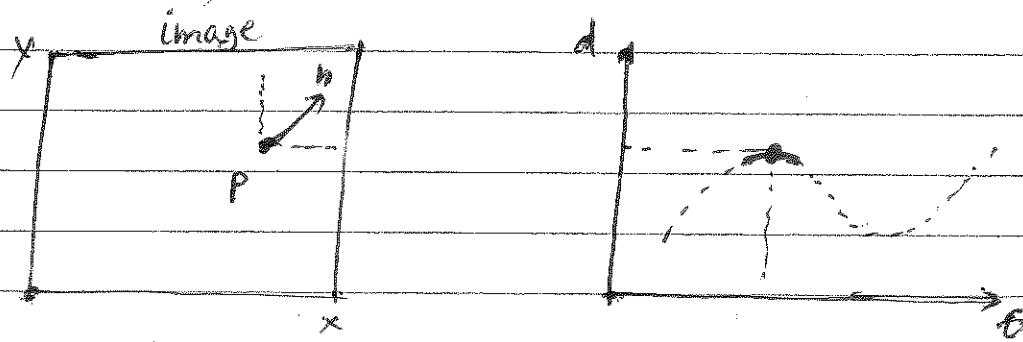
3) Peak detection : threshold & suppress close peaks.

## \* Hough with edge elements

Estimate a normal vector at each pixel



Each pixels provide a vote which is a small curve section (accounting for inaccurate normal direction)



A point  $(x, y)$  with normal  $\theta$  votes for

$$d = x \cos(\theta + \alpha \Delta\theta) + y \sin(\theta + \alpha \Delta\theta)$$

where  $\alpha \in [-1, 1]$

## \* Generalized Hough Transform

→ Objective :  $f(x, y, a_1, \dots, a_m) = 0$

$$\text{eg: } f(x, y, \theta, d) = x \cos \theta + y \sin \theta - d = 0$$

→ for each  $(x, y)$  scan  $a_1, \dots, a_{m-1}$  and vote  
for  $a_m$

→ Problems :

- Sparse space with increased dimensions
- Voting inefficient with increased dimensions.

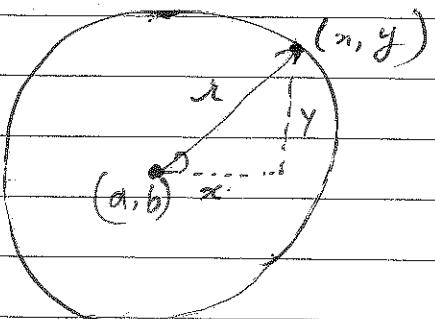
Example for a circle :

$$\rightarrow \text{equation : } (a - x)^2 + (b - y)^2 = r^2$$

or

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

$$y = b + r \sin \theta,$$



→ vote :

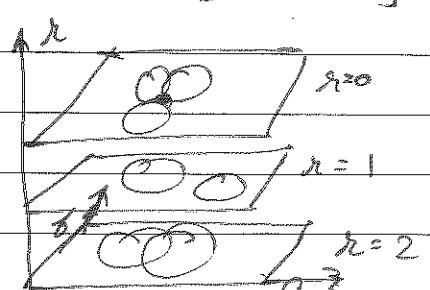
given  $(x, y)$  vote in each  $\theta$  plane using:

$$a = x - r \cos \theta$$

where  $\theta \in [0, 360]$

$$b = y - r \sin \theta$$

~~X-Hough over here X~~



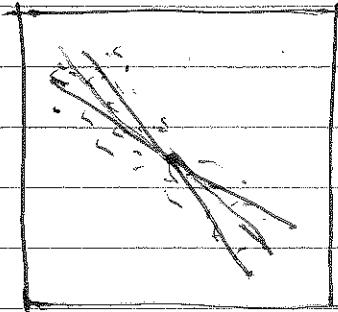
# \* Model fitting

After rough model

Given data  $\{(x_i, y_i)\}_{i=1}^m$

& model  $M_p(x_i)$ :

- ~~determine~~ <sup>define</sup> objective:  $E(p) = \sum_i (M_p(x_i) - y_i)^2$



- minimize objective:  $\hat{p}^* = \underset{p}{\operatorname{arg\,min}} E(p)$

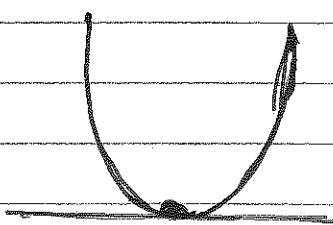
Line fitting: given a group of points (e.g. neighboring a line detected using Hough) find a more accurate model.

given  $\{(x_i, y_i)\}_{i=1}^n$ , find the parameters of line fm

$$y = ax + b$$

$$\left\{ E(a, b) = \sum_{i=1}^n (y_i - (ax_i + b))^2 \right.$$

$$\left. a^*, b^* = \underset{a, b}{\operatorname{arg\,min}} E(a, b) \right.$$



$$\Rightarrow \nabla E(a, b) = 0 \Rightarrow a^*, b^*$$

# Lecture - II Line fitting

Given data  $\{(x_i, y_i)\}_{i=1}^m$  and model  $M_p(x_i)$ :

- define objective :  $E(p) = \sum_i (M_p(x_i) - y_i)^2$

Sum of difference b/w prediction.

- minimize objective :  $p^* = \underset{p}{\operatorname{argmin}} E(p)$

Line fitting :

$$E(a, b) = \sum_{i=1}^n (y_i - (ax_i + b))^2$$

$$\nabla E(a, b) = \begin{bmatrix} \frac{\partial E}{\partial a} \\ \frac{\partial E}{\partial b} \end{bmatrix} = \begin{bmatrix} \sum 2(y_i - ax_i - b)(-x_i) \\ \sum 2(y_i - ax_i - b)(-1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} \sum 2x_i^2 & \sum 2x_i \\ \sum 2x_i & \sum 2 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} \sum 2x_i y_i \\ \sum 2y_i \end{bmatrix}$$

$A \quad X \quad B$

Logic:  $(Ax = b) \times A^{-1}$

$$\cancel{A^{-1}A}x = A^{-1}B$$
$$x = A^{-1}B$$

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\Rightarrow x = A^{-1}b \text{ where}$$

$$A = \begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & m \end{bmatrix} \quad x = \begin{bmatrix} a \\ b \end{bmatrix} \quad b = \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

while doing this in python use = solve(A, B)  
instead of = inv(A).B

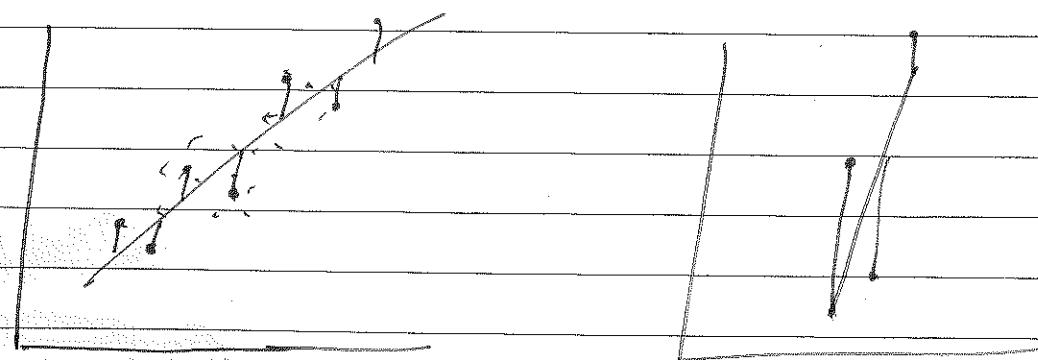
Poor model -

- the solution obtained is optimal in the sense of minimizing the objective:

$$E(a, b) = \sum_{i=1}^n (y_i - (ax_i + b))^2$$

algebraic distance

- The problem is objective itself.
- This take vertical error, which is ~~no~~ bad.



- So this becomes inaccurate.

Solution: To we  $\perp$  projection to a line i.e geometric distance.

Minimize geometric distance:

- line in homogeneous coordinates

- point P on the line l if:

$$\begin{aligned} \ell^T p = 0 \\ (\underline{a}, b, c)^T \quad (x, y, 1) \end{aligned}$$

for p to be on the line:

$$p \cdot n = d$$

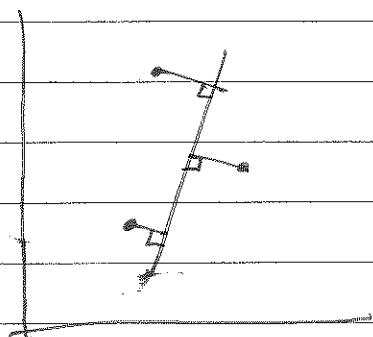
$$nx + by - d = 0$$

\* if  $n = (a, b)$  is not normalized, c is only proportional to distance.

when p is not on the line l  $\ell^T p$  is the geometric distance of p from l

-  $\ell^T p$  is also the algebraic distance  
(deviation from the expected value  
of 0)

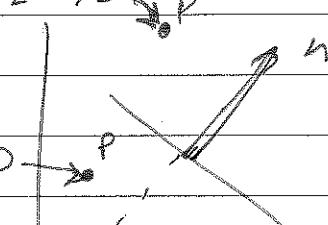
- In line fitting the algebraic error  
has a meaning of being a geometric error.



$$ax + by + c = 0$$

normal  
vector

$\rightarrow$  -ve dist  
from origin



$$LTP < 0$$



$$LTP > 0$$

\* Geometric Line fitting :

$$\{ E(l) = \sum_i^m (l^T p_i)^2$$

$$p_i = (x_i, y_i, 1)$$

$$l^* = \underset{l}{\operatorname{arg\min}} E(l)$$

$$\begin{matrix} 3 \times 3 \\ \downarrow \\ 3 \times 1 & 1 \times 3 \end{matrix}$$

$$E(l) = \sum l^T p_i p_i^T l = l^T \underbrace{\sum p_i p_i^T}_S l = \underbrace{l^T S l}_{1 \times 3 \quad 3 \times 3 \quad 3 \times 1}$$

$$\nabla E(l) = 0 \Rightarrow S l = 0$$

$\Rightarrow l$  = eigenvalues of  $S$  belonging to zero eigenvalue

[eigen value belonging to smallest eigen value.]

To evaluate solution compute  $(E(l))_m$  error/point in pixels

\* correlation matrix (structure tensor)

$S = \sum p_i p_i^T$  = correlation matrix of points  $p_i(x_i, y_i)$

$$S = \begin{bmatrix} \sum x_i^2 & \sum x_i y_i & \sum x_i \\ \sum x_i y_i & \sum y_i^2 & \sum y_i \\ \sum x_i & \sum y_i & m \end{bmatrix} = D^T D$$

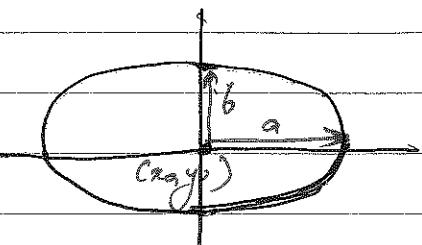
$\Sigma$

$$D = \begin{bmatrix} x_1 & y_1 & 1 \\ x_2 & y_2 & 1 \\ \vdots & \vdots & \vdots \\ x_m & y_m & 1 \end{bmatrix}$$

## Ellipse fitting

Explicit axis-aligned equation:

$$\frac{(x-x_0)^2}{a^2} + \frac{(y-y_0)^2}{b^2} = 1$$

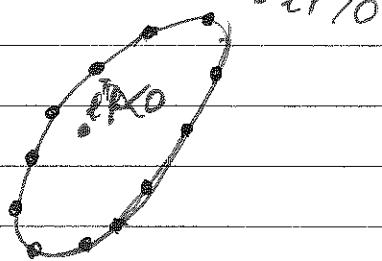


Circle equation is

$$(x-x_0)^2 + (y-y_0)^2 = r^2 \Leftrightarrow \frac{(x-x_0)^2}{r^2} + \frac{(y-y_0)^2}{r^2} = 1$$

\* Implicit equation (conic curve)

$$\left\{ \begin{array}{l} ax^2 + bxy + cy^2 + dx + ey + f = 0 \\ b^2 - 4ac < 0 \end{array} \right.$$



$$\begin{matrix} l^T p = 0 \\ \uparrow \\ (a, b, c, d, e, f) \end{matrix} \quad (x^2, xy, y^2, x, y, 1)$$

convert input points

$$\{(x_i, y_i)\} \rightarrow \{p_i\} \quad p_i = (x^2, xy, y^2, x, y, 1)$$

\* Find parameters  $\ell$  by minimizing algebraic distance:

$$\{ E(\ell) = \sum_{i=1}^m (\ell^T p_i)^2 = \ell^T S \ell$$

$$S = \sum p_i p_i^T \in 6 \times 6$$

$$\ell^* = \underset{\ell}{\operatorname{arg\,min}} \quad E(\ell) \quad \text{such that } b^2 - 4ac < 0$$

\* Rewriting the constraint:

$$b^2 - 4ac < 0 \iff \ell^T C \ell = -1$$

$$C = \begin{bmatrix} 0 & 0 & -2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ -2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

\* Modified objective:

$$\{ E(\ell) = \ell^T S \ell + \lambda (\ell^T C \ell + 1)$$

$$\ell^* = \underset{\ell}{\operatorname{arg\,min}} \quad E(\ell)$$

lagrange multiplier

$$\nabla E(\ell) = 0$$

$$\Rightarrow \lambda S \ell = \lambda C \ell$$

$$\Rightarrow \ell = \lambda S^{-1} C \ell \quad \Rightarrow \ell^* \text{ is the eigenvector of } S^{-1} C \text{ belonging to its -ve eigenvalue}$$

\* algebraic distance of ellipse:

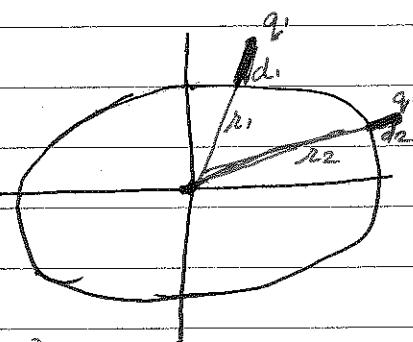
When  $x_i$  is on the ellipse:  $\ell^T x_i = 0$

$\leftarrow 6x_1$

- when  $x_i$  is off the ellipse:  $g_i = \ell^T x_i$   
provides a measure for distance  
from the ellipse.

$$\text{algebraic} = g_i = \ell^T x_i \approx \frac{d_i}{d_i + r_i}$$

distance



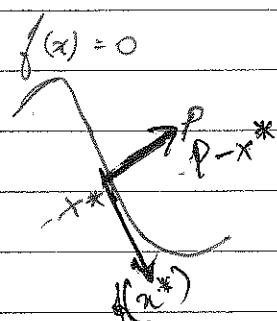
$d_1 \approx d_2$  but  $g_1 > g_2$  (because  $x_2 > x_1$ )

\* geometric distance of ellipse:

geometric distance from a point to arbitrary  
implicit curve:  $f(x) = 0$

$$\{ d(p, f) = |p - x^*|$$

$x^*$  is closest point



\* solution:

$$d(p, f) = |p - x^*| = \frac{|f(p)|}{|\nabla f(x^*)|}$$

algebraic  
distance

To find  $x^*$  solve:

$$\{ \begin{array}{l} f(x^*) = 0 \\ (p - x^*) \cdot \nabla f(x^*) = 0 \end{array}$$

$\leftarrow$  tangent at  $x^*$

\* To find tangent at  $x^*$

1) compute gradient :  $\nabla f(x^*) = \begin{bmatrix} \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \end{bmatrix}^T$

2) Rotate by  $-90^\circ$  :  $t = R(-90) \nabla f(x^*) = \begin{bmatrix} \frac{\partial f}{\partial y}, -\frac{\partial f}{\partial x} \end{bmatrix}^T$

\* Approximated solution :

$$d(p, f) = |p - x^*| = \frac{|f(p)|}{|\nabla f(x^*)|} \approx \frac{|f(p)|}{|\nabla f(p)|}$$

geometric distance                                  algebraic distance

→ reduce the algebraic distance from points with large gradient  $|\nabla f(p)|$

\* Update geometric distance fitting :

$$E(l) = \sum_i \frac{|f(p_i, l)|}{|\nabla f(p_i, l)|} \quad \frac{f(l^T p_i)^2}{\sqrt{(l^T p_i)^2}}$$

where  $f(p_i, l) = l^T p_i$

$$l^* = \underset{l}{\operatorname{arg\min}} E(l) \text{ st. } b^2 - 4ac \leq 0$$

No explicit solution to  $\nabla E(l) = 0$   
 $\Rightarrow$  iterative numerical solution (e.g. gradient descent).

# Lecture 12

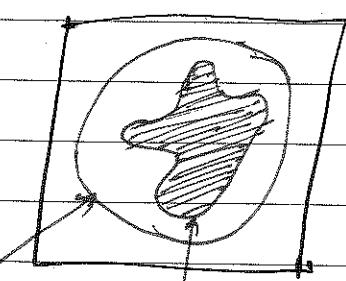
## \* Arbitrary model fitting (eg. splines)

- Given grouppoints  $\{p_i\}$  and a parametric model  $M(p_i, l)$  with parameters  $l$ :
- Define distance between point  $p_i$  & model:
 
$$d_i = d(M(p_i, l), p_i)$$
- Define objective:
 
$$E(l) = \sum d_i^2 = \sum d^2(M(p_i; l), p_i)$$
- Minimize objective to find parameters
 
$$l^* = \underset{l}{\operatorname{argmin}} E(l)$$

## \* Other models

### \* Active Contours (snakes)

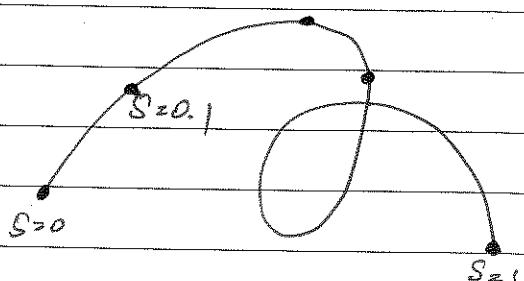
- gradually deform initial contour to fit object boundaries.
- Iterations



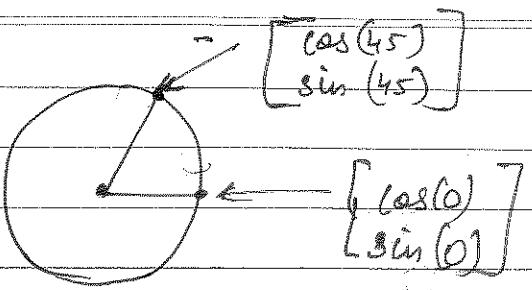
- Use parametric curves  $\phi(s)$  to represent curves
- Initial Contour
- object

Parametric curve

$$\phi(s) = \begin{bmatrix} x(s) \\ y(s) \end{bmatrix}$$



$$\text{Eq: } \phi s = \begin{bmatrix} \cos(s) \\ \sin(s) \end{bmatrix}$$



\* error functional

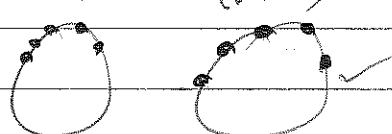
The unknown we are seeking is a function  
 $\phi(c) \Rightarrow$  error measure functional

$$E[\phi(s)] = \int_{\phi(s)} (\alpha(s) E_{\text{cont}} + \beta(s) E_{\text{curr}} + \delta(s) E_{\text{ring}}) ds$$

internal energy
external  
energy

$\alpha(s)$ ,  $\beta(s)$ ,  $\gamma(s)$  are coefficients of the different energy terms.

$$\text{continuity energy: } E_{\text{cont}} = \left| \frac{d\phi}{ds} \right|^2$$



$$\text{Curvature energy : Ecurve} = \frac{d^2\phi}{ds^2} |^2$$

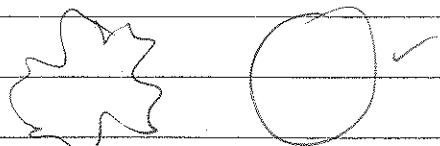



Image energy:  $E_{\text{img}} = -|\nabla I|^2$



\* Discrete functional s-

Discrete case :  $\phi(s) \rightarrow \{P_i\}_{i=1}^n$

$$E_{cont} = \left| \frac{d\phi}{ds} \right|^2 = |P_{i+1} - P_i|^2$$

$$\begin{aligned} E_{curv} &= \left| \frac{d^2\phi}{ds^2} \right|^2 = |(P_{i+1} - P_i) - (P_i - P_{i-1})|^2 \\ &= |P_{i+1} - 2P_i + P_{i-1}|^2 \end{aligned}$$

$$E_{img} = -|\nabla I|^2$$

coordinates unknown  $\rightarrow$  model parameters

Avg distance b/w points to prevent shrinking

$$\begin{aligned} E(\{P_i\}) &= \sum_{i=1}^n \alpha_i (|P_{i+1} - P_i| - d)^2 + \sum_{i=1}^n \beta_i |P_{i+1} - 2P_i + P_{i-1}|^2 \\ &\quad - \sum_{i=1}^n \gamma_i |\nabla I(P_i)|^2 \end{aligned}$$

$\alpha_i, \beta_i, \gamma_i$  are user selected parameters.  $\{P_i\}$  are the unknown model parameters.  $d$  is a computed parameter.

(also zero  $\beta_i$  at corners to allow discontinuity)

Set  $\beta_i$  to zero if curvature @  $P_i$  is  $> \tau$

\* To minimize the error functional  $E$ :

- start with initial guess  $\{P_i\}_{i=1}^n$

- compute coverage distance between points

- for each point  $P_i$  try to minimize  $E$  by testing what happens when moving to neighboring locations

- Repeat ~~step~~ while  $E$  is decreasing

\* Selecting coefficient

- select  $\alpha, B_i, \gamma_i$  to normalize the different terms to a similar scale

- To allow for piecewise curves, set  $B_i = 0$  to points with high curvature:

i.e.

$$\text{when } |P_{i+1} - 2P_i + P_{i-1}| > \tau$$

# A Robust Estimation

## \* Naive approach:

- fit model to all points
- (- compute distance at each point from model
  - discard points with largest distance
  - fit model to remaining points

problems initial model is inaccurate and so we drift in the wrong direction.

## Approaches:

- M-estimates
- RANSAC

## \* M-estimates

- mean square error (mse) fitting:

$$F(\theta) = \sum d^2(x_i; \theta) \quad \text{eg: } d^2 = (e^T x_i)^2$$

for line fitting

- robust estimates:

$$F(\theta) = \sum \beta_a(d(x_i; \theta))$$



MSE is a special case where  $f_a(x) = x^2$

## Geman - McClure estimator

Better than robust

$$g_a(x) = \frac{x^2}{x^2 + \sigma^2}$$

$$x \gg \sigma \Rightarrow g_a(x) = 1$$

$$x \ll \sigma \Rightarrow g_a(x) = \frac{x^2}{\sigma^2}$$

larger  $\sigma$



wider valley

$$E(\theta) = \sum g_a(d(x_i; \theta))$$

$$\nabla E(\theta) = \sum \frac{\partial}{\partial d} \cdot f_g(d) \frac{\partial}{\partial \theta} d(\theta)$$

$$\text{for, } g_a(d) = d^2 \Rightarrow \frac{\partial g}{\partial d} = 2d$$

$$g_a = \frac{d^2}{d^2 + \sigma^2} \Rightarrow \frac{\partial g}{\partial d} = \frac{2d\sigma^2}{(d^2 + \sigma^2)^2}$$

## \* Selecting bandwidth parameter

- large  $\delta \rightarrow$  include more points
- small  $\delta \rightarrow$  include fewer points
- variable estimation:  
start with large  $\sigma$  and decrease as converging.

$$\delta^{(n)} = 1.5 \text{ median } \{ d(x_i; \theta^{(n-1)}) \}$$

↑

estimate at step n

↑

parameters at step n-1

## M-estimator summary:

- 1) Draw a large set of points uniformly at random
- 2) select initial value of  $\theta$
- 3) Fit model  $\rightarrow \theta^{(i)}$
- 4) compute  $\delta^{(i)}$  using median distance of points
- 5) continued while objective is decreasing

\* To overcome incorrect initial guess of  $\theta$  repeat several times and select best solution (smallest objective)

## \* RANSAC:

Random Sample Consensus (RANSAC):

- perform multiple experiments
- choose best results
- use small sets in hope that at least one set will not have outliers.

Parameters:

$n$  = # points drawn at each evaluation

$d$  = min # points needed to estimate model

$K$  = # trials

$t$  = distance threshold to identify inliers

## \* Repeat $K$ times:

- Draw  $n$  points uniformly at random (with replacement)
- fit a model to points
- Find inliers in entire set [distance  $\leq t$ ]
- Recompute model [if at least  $d$  inliers]
- update parameters  $(K, t)$

## \* Choose best solution:

→ largest consensus set

- (as smallest error)

# Lecture 13 - 10711

(Model fitting (contd) - - - )

\* How to identify threshold ( $t$ ) :

to estimate  $t$  use median distance from model.

$$t = 1.5 \times \text{median}$$

\* To estimate ( $K$ ) use :

(user)  $P$  : With probability of  $p$  at least one (selected) experiment does not have outliers (e.g.  $p=0.99$ )

(estimated)  $\omega$  : probability that a point is an outlier  
(initially  $\omega = 0.5$ )

So probability that all  $k$  experiments failed :

$$(1-P) = (1-\omega^n)^k$$

$$\log(1-P) = k \log(1-\omega^n)$$

$$k = \frac{\log(1-P)}{\log(1-\omega^n)}$$

large  $P \rightarrow$  large  $k$   
small  $\omega \rightarrow$  large  $k$

$$\omega \leftarrow \frac{\# \text{ outliers}}{\# \text{ points}}$$

update  $\omega, K$  every iteration  
but set upper bound for  $K$ .

# Segmentation and

## Classification

A  
Segmentation :

2 class segmentation - foreground & Bg

3 class segmentation - detect vegetables from 3 colors.

It is difficult to do with just colors.

Semantic segmentation - most difficult - identify everything as class.

Perceptual segmentation problem (more difficult) :  
Computer cannot perceptualize objects.

\* Problem Statement :

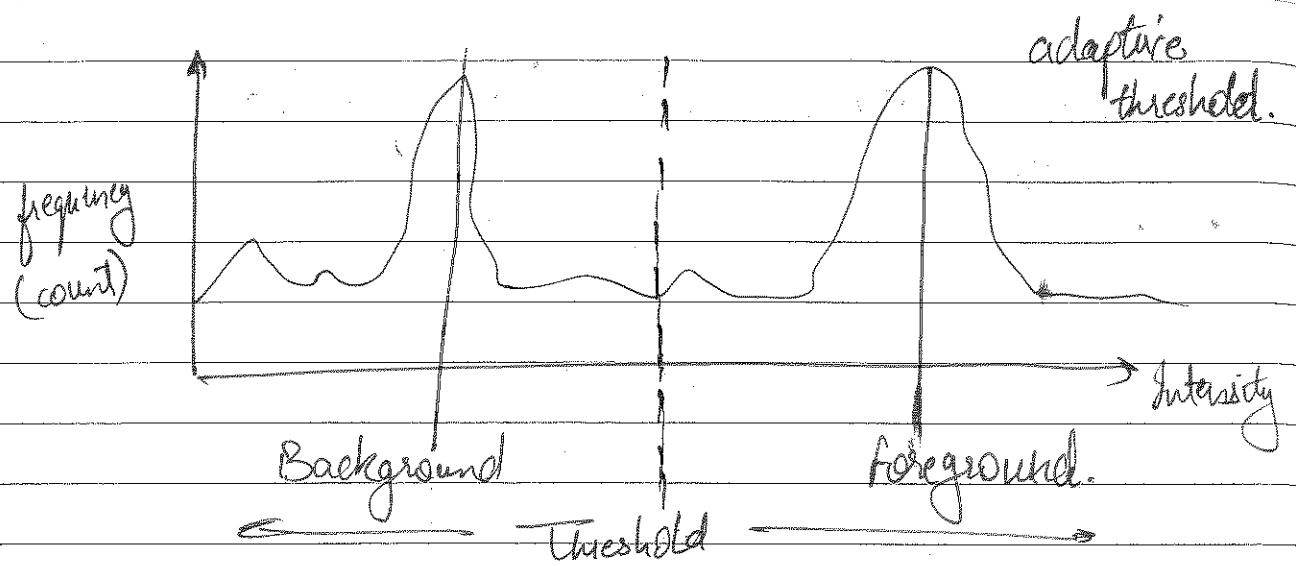
- Separate objects from background.
- Find contours of object
- Semantic image segmentation : label each pixel in the image with class label.

\* Color segmentation

- Use color thresholds to identify objects

Problem : how to identify color threshold

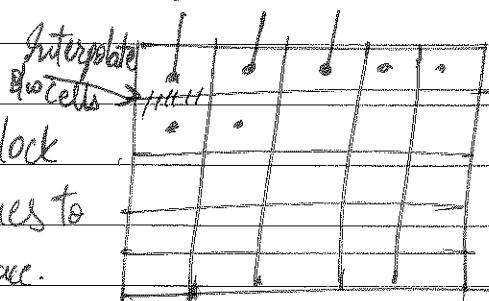
Solution : Use Intensity histogram.



\* But this can still skip certain Objects,

Solution : use cells

- Break image into cells
- find threshold in each block
- Interpolate threshold values to create segmentation surface.
- at each pixel location compare image intensity to segmentation surface.



So this is adaptive threshold.

\* Special context : to group pixels together which have similar characteristics. This gives super pixels.

Agglomerative segmentation :

Start with each pixel in a separate cluster

- Merge clusters with small distance
- color & char difference
- Repeat while clusters are not satisfied

Super clusters

satisfactory

### Divisive segmentation:

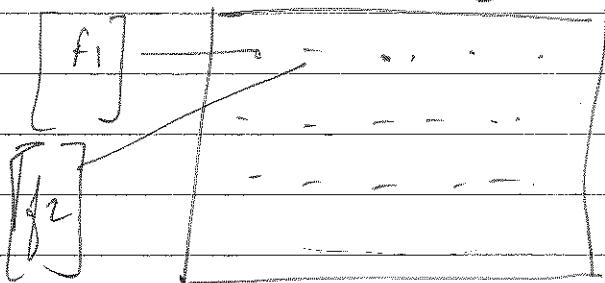
- start with all pixels in one cluster
- split ~~the~~ clusters to provide large distance between them.
- Repeat while clusters are not satisfactory.

### \* Feature based segmentation

- Define feature vector at each pixel  $x$ :

$$f(x) = \begin{bmatrix} x \\ I(x) \\ L(x) \end{bmatrix} \quad \begin{array}{l} \xleftarrow{\hspace{1cm}} \text{location} \\ \xleftarrow{\hspace{1cm}} \text{intensity} \\ \xleftarrow{\hspace{1cm}} \text{local characteristic (eg texture)} \end{array}$$

- Apply clustering



### → Clustering (Machine Learning Concept)

- each pixel is assigned a pixel feature vector
- To define a cluster group pixels with similar feature vectors.
- The similarity ~~is~~ in each cluster should be high compared to the similarity between clusters.

## \* fundamental problem in clustering :

- Need to compute cluster parameters and assign items ~~and~~ to clusters
- Joint estimates is difficult.

## \* Issues :

- what is the number of clusters?
- How to compute similarity within clusters and distance between clusters?
- what features to use?

## \* K means (clustering algo) :

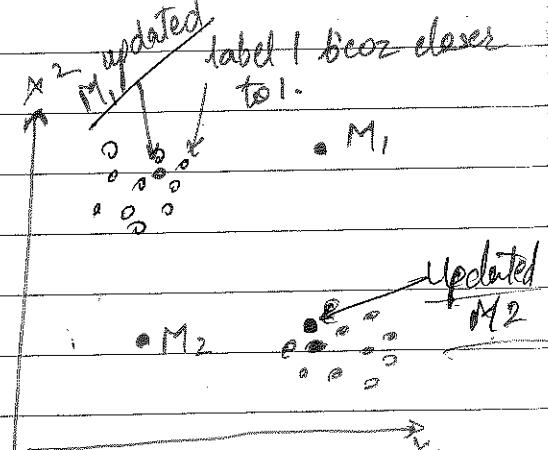
- select K
- select initial guess of means :  $m_1, \dots, m_k$
- Repeat while  $m_j$  change:

$$\textcircled{1} \quad l_i = \underset{j \in [1, K]}{\operatorname{argmin}} \|f_i - m_j\|^2$$

for each price assign labels

$$\textcircled{2} \quad S_j = \{i \mid l_i = j\}$$

$$\textcircled{3} \quad m_j^* = \frac{\sum_{i \in S_j} f_i}{\# S_j} \quad \leftarrow \begin{array}{l} \text{average new center} \\ (\text{compute mean}) \end{array}$$



\* simultaneous solution of the two problems :

- find cluster centers
- assign points to clusters

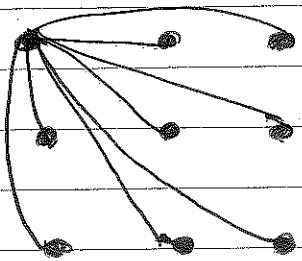
hold one parameter fixed and change other to minimize error.

$$\Rightarrow \text{coordinate descent on } E(\{\mathbf{c}_{ij}\}, \{\mathbf{m}_j\}) = \sum \|f_i - m_j\|^2$$

\* graph cuts (other method than clustering) :

to identify clusters / super pixels.

- Represents image using graph.
- each pixel is connected to all other pixels.
- The weight on the link between nodes  $p$  &  $q$  represents the similarity b/w them (intensity + location).



$$w_{pq} = \exp(-\|f(p) - f(q)\|)$$

if  $f(p)$  similar to  $f(q)$  then  $w_{pq} = 1$   
else  $w_{pq} = 0$

- Delete links to create segments.
- Remove low similarity links.
- Nodes in each cluster should have strong links.



\* possible graph cuts remove links to create disconnected subgraphs.

- The cost of a cut is the sum links being removed:

$$\text{cut}(A, B) = \sum_{p \in A, q \in B} w_{p,q}$$

- Minimum cut: cut with lowest cost.

\* But this is not desirable as each pixel will end up being 1 subgraph.

- a normalized cut assigns cost taking into account the size of produced clusters.

$$\text{Ncut}(A, B) = \frac{\text{cut}(A, B)}{\text{vol}(A)} + \frac{\text{cut}(A, B)}{\text{vol}(B)}$$

$$\text{vol}(A) = \sum_{p \in A, p \in A \cup B} w_{p,q}$$

\* finding normalized cuts:

- the weighted degree  $d_i$  of a node is the sum of all edges connected to it:

$$d_i = \sum_{j=1}^n w_{ij}$$

weighted degree of each node  $\rightarrow D = [d_1, \dots, d_n]$

~~Similarity matrix~~  $\omega = [w_{11}, \dots, w_{1n}; \dots; w_{n1}, \dots, w_{nn}]$

= Defining a cut :

$$X = \{1, -1\}^n$$

$$x(i) = 1 \iff i \in A$$

$$x(i) = -1 \iff i \in B$$

$$Ncut(A, B) = \frac{\text{cut}(A, B)}{\text{vol}(A)} + \frac{\text{cut}(A, B)}{\text{vol}(B)}$$

$$= \underbrace{\sum_{\substack{x_i > 0 \\ x_j < 0}} (-w_{ij} x_i x_j)}_{\sum_{x_i > 0} d_i} + \underbrace{\sum_{\substack{x_i > 0 \\ x_j < 0}} (-w_{ij} x_i x_j)}_{\sum_{x_j < 0} d_j}$$

\* Define :

$$K = \frac{\sum_{x_i > 0} d_i}{\sum_i d_i} \quad b = \frac{K}{1-K} = \frac{\sum_{x_i > 0} d_i}{\sum_{x_i < 0} d_i}$$

$$y = (1+x) - b(1-x)$$

Convert the problem :

$$\min_x Ncut(x) \approx \min_y y^T (D-w) y$$

$$\text{s.t } y^T D y = 0$$

$$1 \equiv [1 \dots 1]$$

\* since it's hard to solve with discrete as  $[1, -1]$  needs to be assigned, we make it continuous domain.

$$\text{Solve: } \min_{\mathbf{y}} \frac{\mathbf{y}^T (\mathbf{D} - \omega) \mathbf{y}}{\mathbf{y}^T \mathbf{D} \mathbf{y}} \quad \text{s.t. } \mathbf{y}^T \mathbf{D} \mathbf{1} = 0$$

highlight gradient

$$\Rightarrow (\mathbf{D} - \omega) \mathbf{y} = 2 \mathbf{D} \mathbf{y}$$

- the first eigen vector (belonging to  $\lambda = 0$ ) is  $\mathbf{1} = [1 \ -1 \ \dots \ 1]^T$

- the second smallest eigenvector is the solution.  
(the one that is not zero)

## \* Spectral Clustering (like graph cut somewhat):

- has a degree matrix as before.
- has an adjacency matrix ( $\omega$ ) which is very similar but just that the similarity matrix have a value of 1 else a 0.
- graph laplacian:  $L = D - \omega$
- $L$  is symmetric and positive semidefinite
- the sum of every row & column in  $L$  is  $0$ .
- the e-vector belonging to second eigen value as smallest value (0) has a vector  $[1 \dots 1]$
- For every  $x \in \mathbb{R}^n$ :  $x^T L x = \sum_{i=1}^n \sum_{j=1}^n \omega_{i,j} |x_i - x_j|^2$

## \* So to compute spectral cluster:

- compute eigen vector of  $L$
- Stack  $K$  smallest eigen vectors:  $U = [u_1 \ u_2 \ \dots \ u_K]$   
with non zero eigen values.  
 $n \times K$   
 $K$  dimension
- define  $y_i$  as the  $i$ -th row of  $U$   
(i.e. map data to  $K$ -dimensional space)  $n \times D \rightarrow K \times D$   
cluster the  $n$  rows  $y_i$  with  $K$ -means  
(i.e. each cluster in  $K$ -D space instead of  $n$ -D)

## \* Classification:

- Classification
- classification + localization
- object detection
- Instance segmentation
- semantic segmentation

## \* Problem with object recognition?

- Semantic gap (pixels vs labels):
  - pixels ~~with~~ values as descriptors are sensitive to small variations.
- Object ~~recognition~~ needs to be invariant
  - pose
  - illumination
  - deformation (articulated objects)
  - occlusion
  - background
  - natural variability (e.g. cars)

## \* Want:

- generalization (invariance)
- extend to other problems

\* Need:

- feature extraction (invariance)
- classification algorithm.

→ Some traditional object recognition methods:-

today state of the art is to use deep learning

Eigen faces:

- map image (eg.  $100 \times 100$ ) to lower dimensional vectors (eg 64) using PCA.
- measure similarity to templates in lower dimensional space (less sensitive to small variations.)

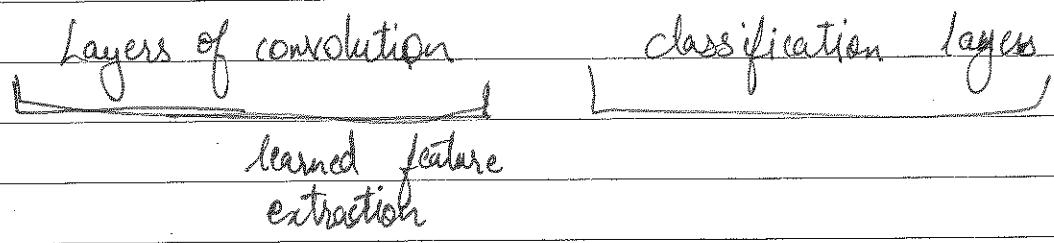
Bag - of - words

- extract features from image (eg SIFT or HOG)
- cluster features to create a codebook (dictionaries)  
compute a distribution of code words in each class
- classify using distribution of code words.

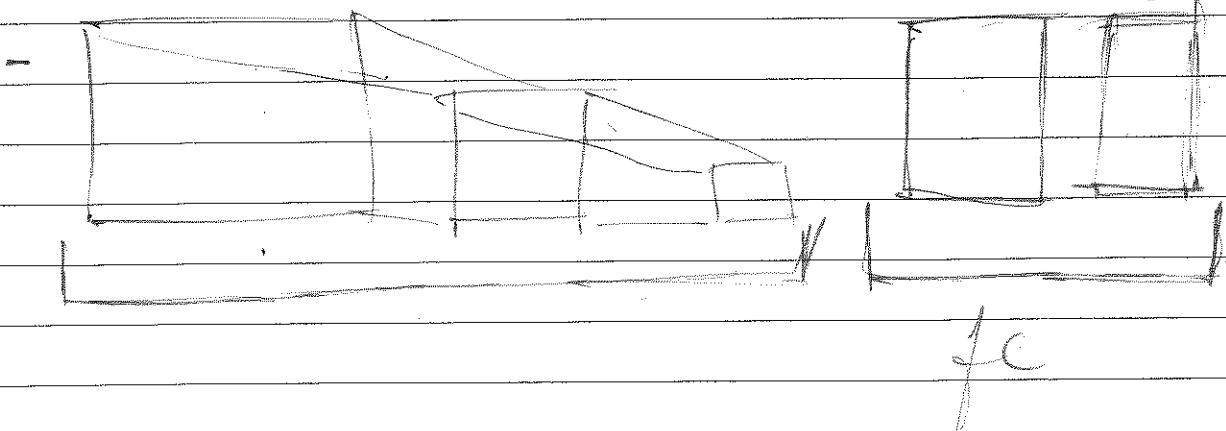
## \*Convolutional neural networks (CNN):

codebook Convolutional neural networks (CNN) learn feature extraction and classification

- learn from examples instead of coding algorithms
- Architecture for one class can work for others by changing data  $\rightarrow$  class generalization.



- Combine features at multiple scales (multiple scale analysis)



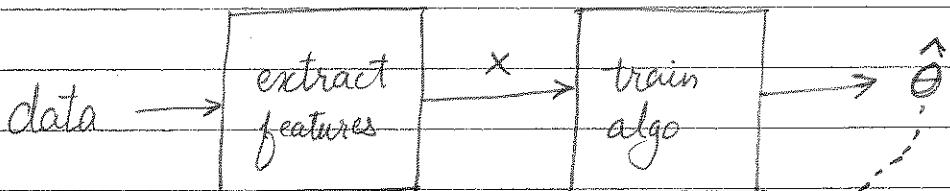
# Lecture 14



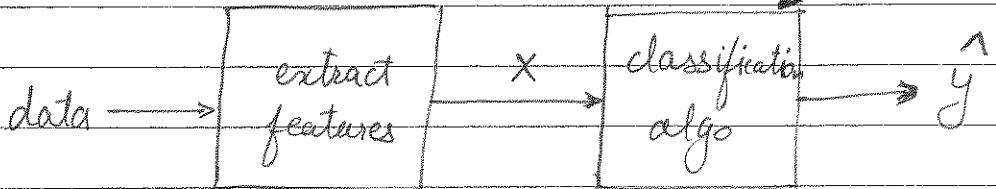
## DEEP LEARNING (DATA-DRIVEN)

\* Training & inference :-

Training :-

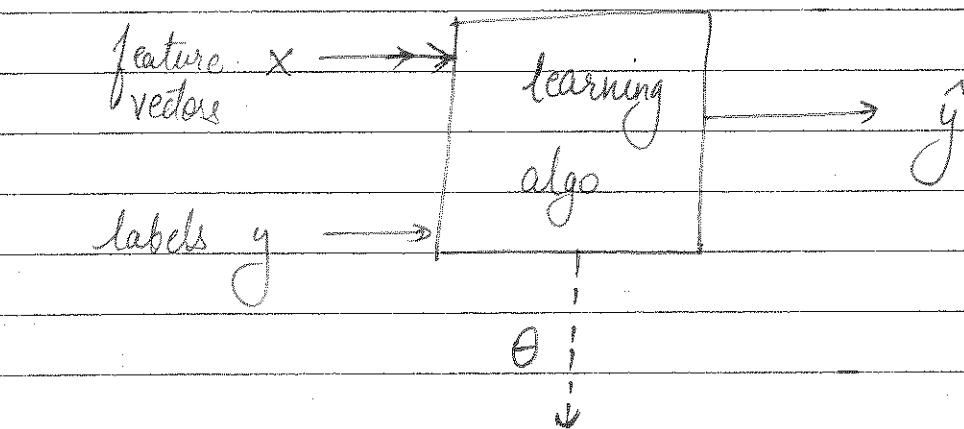


Inference :-



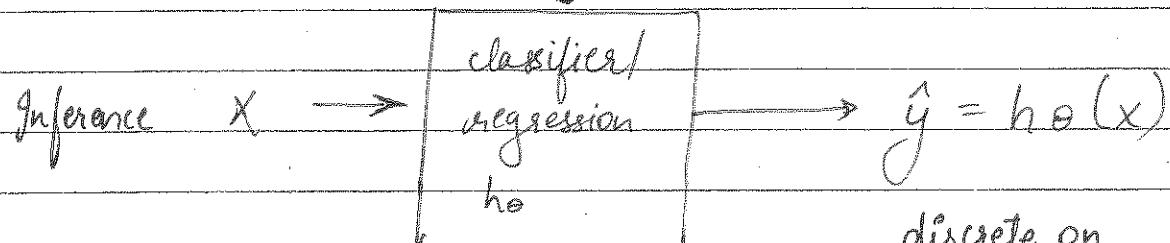
\* Supervised learning (common method) :-

$$\{x^{(i)}, y^{(i)}\}_{i=1}^m \in \mathbb{R}^n$$



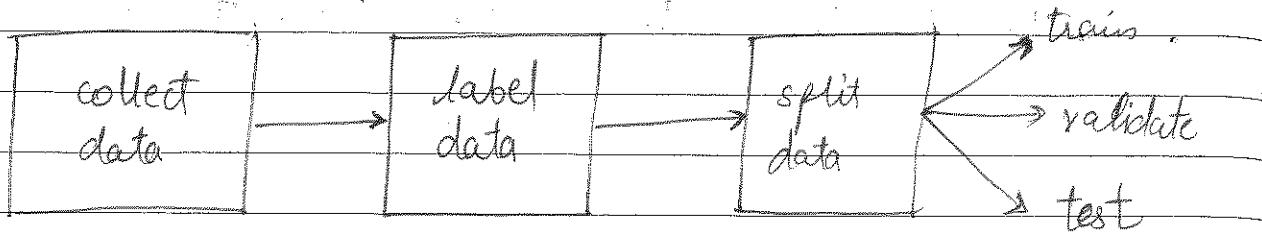
$$\theta = \underset{\theta}{\operatorname{arg\! min}} L(\theta)$$

loss function  
compare  $y$ ,  $\hat{y}$



discrete or  
continuous output

## \* Data collection :



- Train / test on different collection (eg. K fold cross validation for small data available)
- Performance based on testing error (generalization vs. overfitting)
- Large models have more capacity but require more data (exponentially) and may overfit.
- Parameters :
  - Learn model parameter (training)
  - Hyper-parameters are specified parameters.

## Deep Learning

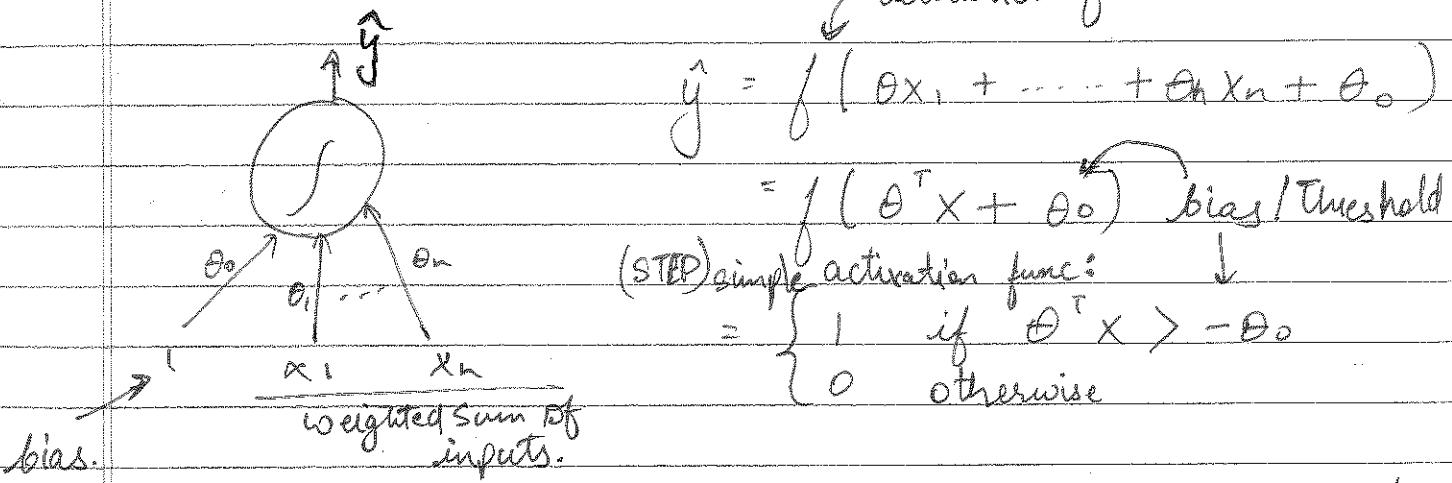
- Alexnet
- ResNet
- Google Net / Inception
- Yolo - realtin
- U- net (image segmentation)
- DeepLab (semantic segmentation)
- triple net
- Deep Pose
- VGG (classification - feature extraction)
- Cycle GAN (Image mapping)

→ Deep dream (Image hallucination)

→ Image enhancement (chen et al 2018)

## \* Neural Networks

- A neural network is a collection of artificial neurons.
- artificial neurons are linear classifiers with activation



- fire if weighted sum of inputs is greater than threshold/bias.

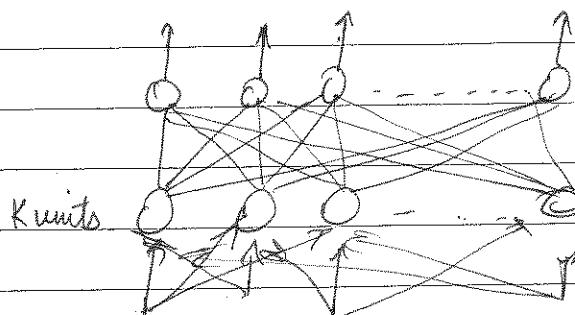
→  $\theta_0, \theta_1, \theta_2, \dots, \theta_n$  are learnt during process

→ They are called linear classifiers, although activation function may not be linear.

$$\hat{y} = f(x; \theta) = \underbrace{\theta^T x}_{K \times 1} + \underbrace{\theta_0}_{K \times 1}$$

Neurons are in layers.

$$\begin{array}{l} \hat{y} > 0 \rightarrow c_1 \\ \vdots > 0 \rightarrow c_2 \\ \vdots > 0 \rightarrow c_3 \\ \vdots \\ g_k : : : \end{array}$$



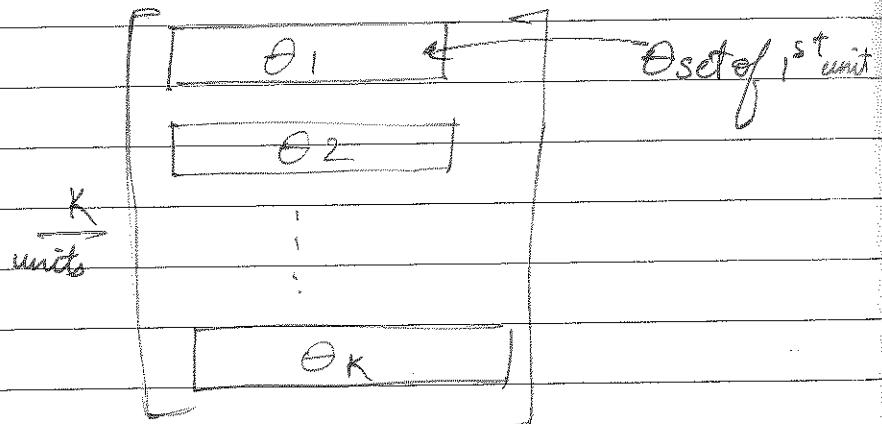
So,

$$\hat{y} = f(x; \theta) = f(\theta^T x + \theta_0)$$

$\theta$   
 $\theta_1$   
 $\theta_2$   
 $\vdots$   
 $\theta_K$

$x$   
 $x_1$   
 $x_2$   
 $x_3$   
 $\vdots$   
 $x_n$

$\theta$  matrix



N no of inputs

so take product,

$$\begin{matrix}
 & \boxed{\theta \text{ matrix}} & \boxed{x} \\
 K & & \\
 N & & \\
 & \text{Input vector} &
 \end{matrix}$$

so far an entire layer we have,

$$\text{layer: } \theta_{K \times N} \text{ Do } \hat{y}_{K \times 1} \text{ } y_{K \times 1}$$

### \* Linear Classifiers: interpretation

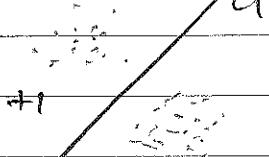
→ perform linear discrimination

each unit defines a linear decision boundary

$K$  planes as

→ perform template matching

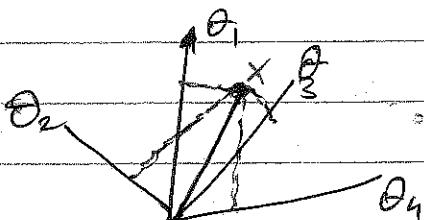
each unit defines a template



→ perform projection into a new basis

each unit represents a basis vector

$\theta$  parameters of plane separating classes. (one for each unit)



- Note: in the end, it's taking a weighted sum and passing it to activation function.

- each linear discriminant separates one class from all others.

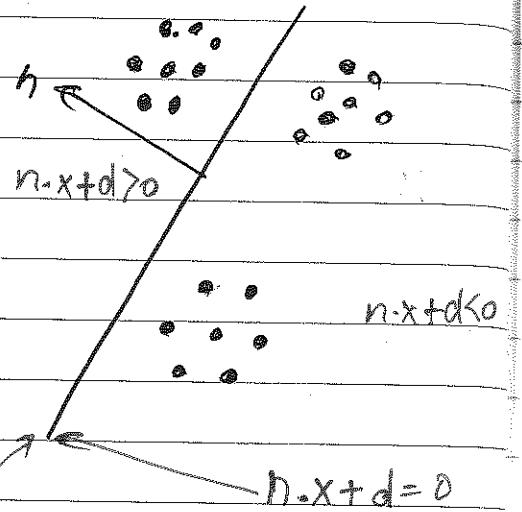
- The value of a linear discriminant is +ve for examples belonging to class and negative otherwise.

- It measures distance from the decision boundary.

\* just like tough; we have a discriminant function,

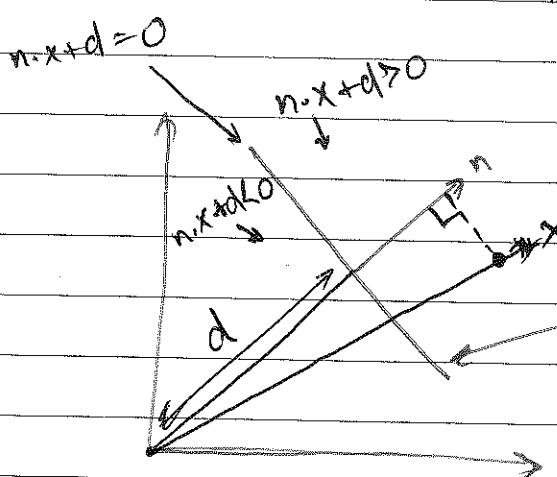
$$\hat{y} = \theta^T x + \theta_0$$

$$= \begin{bmatrix} n_1^T \\ n_2^T \\ \vdots \\ n_k^T \end{bmatrix} \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix} + \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_k \end{bmatrix}$$



$$= \begin{bmatrix} n_1 \cdot x \\ n_2 \cdot x \\ \vdots \\ n_k \cdot x \end{bmatrix} + \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_k \end{bmatrix}$$

hyper plane       $n = \text{normal}$   
 $d \sim \text{negative distance from origin}$



discriminant/  
Hyperplane

### \* Template matching

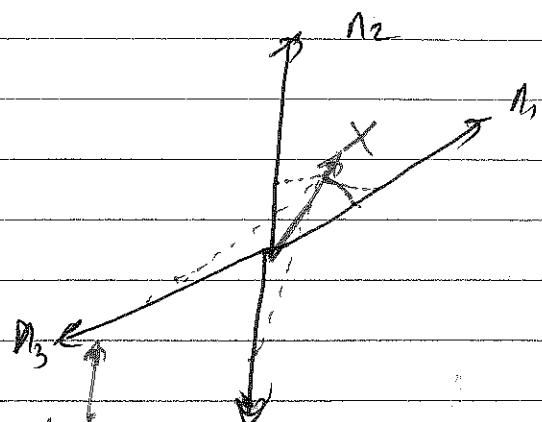
- rows of  $\theta^T$  (rows of  $\theta$ ) are templates.
- with  $K$  rows of  $\theta^T$  we have  $K$  templates (one template per class).
- $\theta^T x$  measures how well  $x$  matches (dot product of  $x$  with rows of  $\theta^T$ ).
- High similarity to a template of a particular class

indicates high membership in this class.

$$\hat{y} = \theta^T x + \theta_0$$

$$= \begin{bmatrix} n_1 \\ \vdots \\ n_k \end{bmatrix} \begin{bmatrix} x \end{bmatrix} + \begin{bmatrix} d_1 \\ \vdots \\ d_k \end{bmatrix}$$

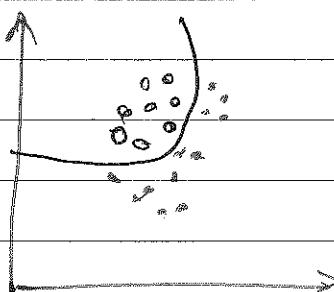
$$= \begin{bmatrix} n_1 \cdot x \\ \vdots \\ n_k \cdot x \end{bmatrix} + \begin{bmatrix} d_1 \\ \vdots \\ d_k \end{bmatrix}$$



\* What if decision boundaries are non linear?  
What if more than one template is needed per class?

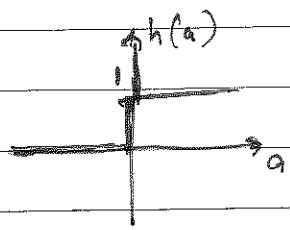
→ combine multiple linear classifiers.

It's important to put activation function to these layers.

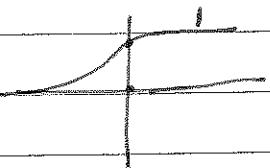


## \* Activation functions

- Step :  $h(a) = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$

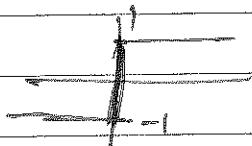


- Sigmoid :  $h(a) = \frac{1}{1 + \exp(-a)}$

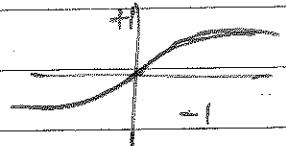


~~adv~~ it does not make sharp decision when its not so certain.

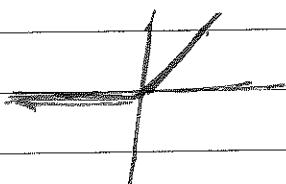
- Step\* :  $h(a) = \begin{cases} 1 & \text{if } a > 0 \\ -1 & \text{otherwise} \end{cases}$



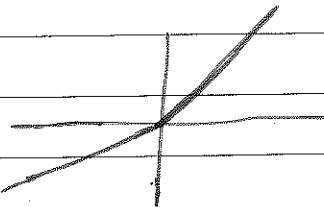
- tanh :  $h(a) = \tanh(a) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



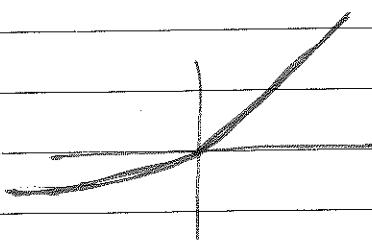
good choice  $\rightarrow$  Relu :  $h(a) = \max(0, a)$   
(Rectified linear unit)



- Leaky Relu :  $h(a) = \max(\alpha a, a)$



- ELU :  $h(a) = \begin{cases} a & a \geq 0 \\ \alpha(e^a - 1) & a < 0 \end{cases}$



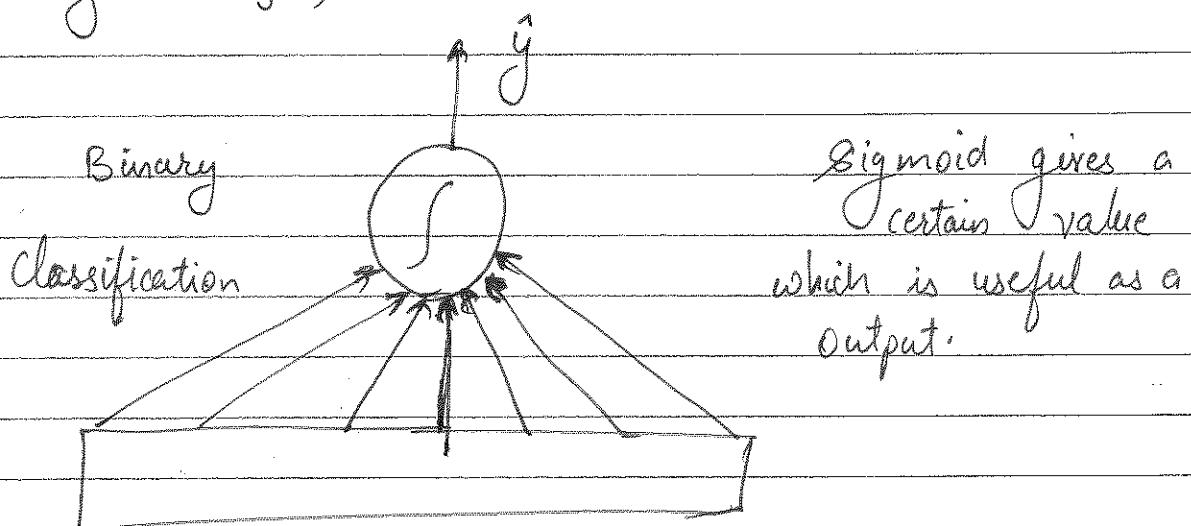
\* Converting similarity to probability:

- Given similarity score (or decision boundary distance  $s_j$ ) from a linear classifier (higher better), compute:

$$\hat{y}_j^{(i)} \equiv p(y=j | x^{(i)})$$

ith example  
label input  
 $\downarrow$   
 $j^{\text{th}}$  unit

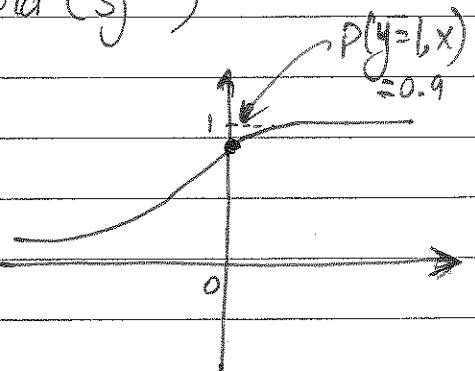
- For the last layer of the network use sigmoid ( $s_j^{(i)}$ ).



- In a 2 class classification case:

$$\hat{y}_j^{(i)} \equiv p(y=j | x^{(i)}) = \text{sigmoid}(s_j^{(i)})$$

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



→ what happens in a k-class classification?

- Use softmax for K-class classification:

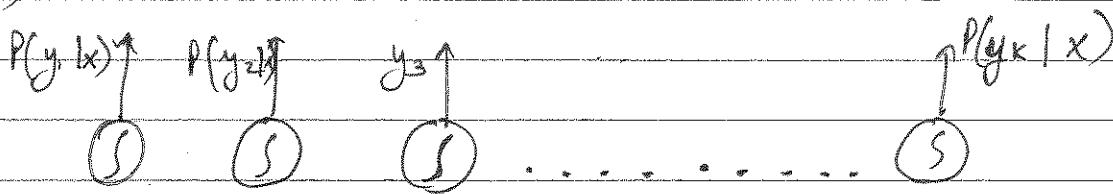
$$\hat{y}_j^{(i)} \equiv P(y=j | x^{(i)}) = \frac{\exp(s_j^{(i)})}{\sum_{l=1}^k \exp(s_l^{(i)})} \leftarrow \text{Softmax}$$

# Lecture 15 - W8L1

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

(contd)



Since probability by taking the sum of their output would be  $> 1$  we need to normalize each outcome. So we take

$$\frac{P(y=j|x)}{\sum P(y=j|x)}$$

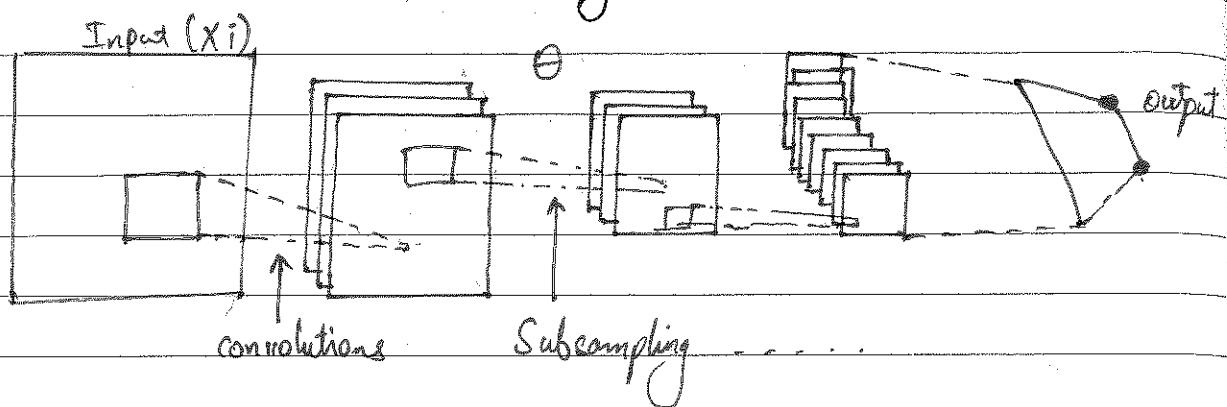
Our outcome.

Softmax

$$y_j^{(i)} = P(y=j | x^{(i)}) = \frac{\exp(s_j^{(i)})}{\sum_{l=1}^n \exp(s_l^{(i)})}$$

(only used at output layer.)

## \* Network Architecture training :



training : find weights  $\theta$  across layers.

input in image ( $x_i$ )

minimize difference between expected & obtained  
output:

$\&$  input with labels  $\{x_i\}, \{y_i\} \quad i=1^m$

- Loss function  $L(\theta)$  is difference b/w outcomes  $y(i) \leftrightarrow \hat{y}(i)$
- Optimization framework is used to minimize loss. By 'how to update properly'.

Gaurav

3

## Loss function

- The weights  $\theta$  are set to minimize loss:

$$\theta^* \equiv \underset{\theta}{\operatorname{argmin}} L(\theta)$$

Sample loss

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_i(x^{(i)}, y^{(i)}; \theta)$$

$$= \underset{\theta}{\text{argmin}} \sum_{i=1}^m L_i(f(x^{(i)}; \theta), y^{(i)})$$

not exactly required

$y^{(i)}$  compare

not exactly required

- L1 and L2 loss:

- In regression problems minimize problems  $\rightarrow$  distance between known and predicted values:

$$L_1 : L_1(\theta) = \sum_{j=1}^k | \hat{y}_j^{(i)} - y_j^{(i)} |$$

X	y	y
-	3.2	-3.1
-	5.3	-5.6
-	7.2	-9
-	;	;
;	;	;

$$L2: L_i(\theta) = \sum_{j=1}^k (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

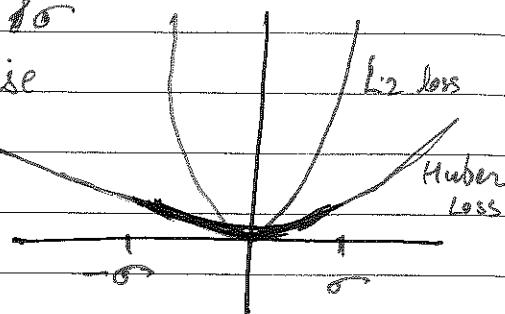
$\hat{y}_j^{(i)}$  is predicted value :  $y_j^{(i)} \in \mathbb{R}$

## Huber Loss (like Geman - McClure)

- Cap loss of outliers (robust regression) and also its easy to take derivative of this function

$$S_\sigma(d) = \begin{cases} \frac{1}{2}d^2 & \text{if } |d| \leq \delta_\sigma \\ \sigma(d - \frac{1}{2}\sigma) & \text{otherwise} \end{cases}$$

- quadratic for small  $d$
- linear for large  $d$



$$L_i(\theta) = \sum_{j=1}^K S_\sigma(\hat{y}_j^{(i)} - y_j^{(i)})$$

## Cross entropy Loss (classification)

\* convert similarity scores to probabilities:

$$\hat{y}_j^{(i)} = P(y=j | x^{(i)}) \quad j \in [1, K]$$

Previous ones were for continuous outcomes, what if we want discrete outcomes (0, 1).

\* Likelihood:  $L(\theta) = \prod_{i=1}^m \prod_{j=1}^K (P(y=j | x^{(i)}))^{y_{j(i)}}$

↑  
maximize      ↑  
                   $y_{j(i)}$   
                   $\hat{y}_j^{(i)}$

(multiply pred. prob. for correct class.)

• Only the relevant  $\hat{y}_j^{(i)}$  outcomes are 1 and rest all become  $x^{(i)}_j = 1$ , so they do not contribute.

\* Problems with likelihood : ① We want to minimize loss & to be consistent with other losses. So we minimize -ve of likelihood. ② It is difficult to take derivative of product so we ~~use~~ need to change it to sum for which we use log.

$$\text{* Log likelihood : } l(\theta) = -\log L(\theta) = - \sum_{i=1}^m \sum_{j=1}^K y_j^{(i)} \log (P(y=j|x^{(i)}))$$

(-ve LL)

minimize

$$l(\theta) = - \sum_{i=1}^m \sum_{j=1}^K y_j^{(i)} \log (\hat{y}_j^{(i)})$$

cross entropy loss

$$\text{where sample loss} = L_i(\theta) = - \sum_{j=1}^K y_j^{(i)} \log (\hat{y}_j^{(i)})$$

for single example

\* possible cross-entropy loss values :

$$L_i(\theta) \in [0, \infty]$$

0 when  $P \rightarrow 1$

$\infty$  when  $P \rightarrow 0$

- For random class assignment (worst case) :

$$P(y=j|x^{(i)}) = \frac{1}{K} \rightarrow -\log (P(y=j|x^{(i)})) = K$$

$$L(\theta) = \log(K)$$

worst loss value

## Regularization

- Occam's razor: simpler explanations are better
  - i.e. generalization because more complex ones can be overfitted & it's not desirable.
- A simpler solution is when the weights  $\theta$  are lower (eg when  $\theta_{ij} = 0$  we remove one coefficient)
- smaller coefficients  $\rightarrow$  more stable solution that will generalize better

$$L(\theta) = \underbrace{\frac{1}{m} \sum_{i=1}^m L_i(f(x_i; \omega), y_i)}_{\text{error term}} + \lambda \underbrace{R(\theta)}_{\substack{\text{weight of term} \\ \text{regularization}}} + \underbrace{\lambda}_{\text{regularization (hyper parameter)}}$$

where  $R(\theta) = \sum_i \sum_j \theta_{ij}^2$

$\uparrow$   $\uparrow$   $\uparrow$

L<sub>2</sub> regularization

$$L_1 \text{ regularization: } R(\theta) = \sum_i \sum_j |\theta_{ij}|$$

II) Now that we have  $L(\theta)$  loss func value, we need to use optimization framework.

## Optimization

- To minimize loss (objective) solve:

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta) \Rightarrow \nabla L(\theta) = 0 \Rightarrow \theta^*$$

- If  $\nabla L(\theta)$  is not linear it is hard to find an explicit solution. (The activation functions are non linear)
- So we use numerical iterative solution (e.g. gradient descent).

## Gradient Descent

- we guess the parameters, compute loss, then make second guess that's ~~better~~ based on first loss gradient.

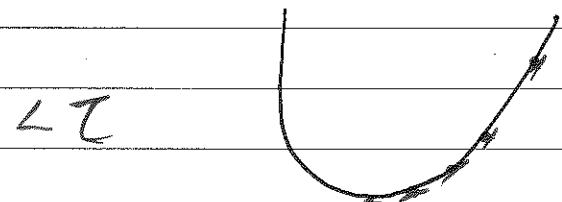
- start with initial guess  $\theta_0$ .

- repeat:

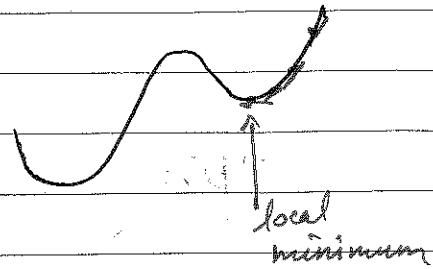
$$\theta^{(i+1)} \leftarrow \theta^{(i)} + -\eta \nabla L(\theta^{(i)})$$

- stop when

$$(L(\theta^{(i+1)}) - L(\theta^{(i)})) < \epsilon$$



we might get local minima and get stuck. So solution to these are available. (eg momentum.)



- The learning rate has to be selected correctly: if too big it overshoots, if its small barely moves but for good it's fine.
- we cannot apply gradient descent to piecewise constant functions:



### \*Example:

- linear regression

$$\hat{y} = \theta^T x$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}_{m \times n}$$

$$L(\theta) = \frac{1}{2} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$\hat{y} = \begin{bmatrix} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(k)} \end{bmatrix}_{m \times 1}$$

$$= \frac{1}{2} (x\theta - y)^T \underbrace{(x\theta - y)}_{m \times 1}$$

$$\nabla L(\theta) = \frac{\partial}{\partial \theta} \frac{1}{2} x^T (x\theta - y)$$

- gradient descent will be:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \begin{matrix} x^T (\hat{y} - y) \\ \text{mxm} \end{matrix}$$

$n \times 1$       scalar

$\hat{y}$   
 $m \times 1$   
 $m \times 1$

- Summation form

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \sum_{j=1}^m \underbrace{\left( x^{(j)T} \theta - y^{(j)} \right)}_{\text{scalar}} x^{(j)}$$

- update parameters by adding examples  $x^{(j)}$   
 where the prediction is wrong.

# Lecture -16 W8L2

## Stochastic Gradient Descent (SGD)

- Randomly order examples
- For  $j = 1 \dots m$

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L_j(\theta^{(i)})$$

gradient based on example  $j$

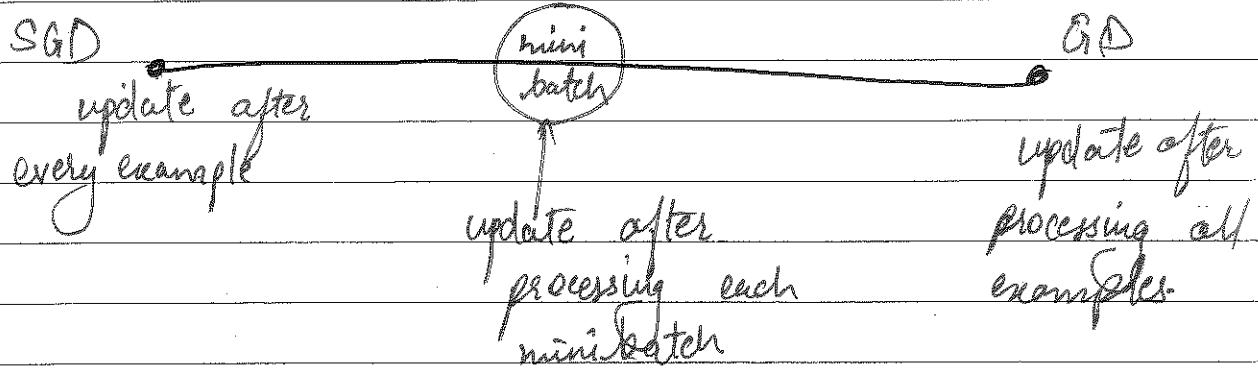
No sum

epoch  $\rightarrow$  traverse entire dataset

This converges faster.

disadvantage is the gradient could be bad or wrong directional which is not desired.

Thus in real world we never perform true SGD, but a ~~hybrid~~ hybrid is used.



## Learning Rate Decay

- Learning rate need not be fixed
- make learning smaller as iterations progress

\* Strategies:

i) Step decay:

every  $K$  iterations  $\eta \leftarrow \eta / 2$

2) exponential decay  $\rightarrow$  decay rate

$$\eta_t = \eta_0 e^{-k/t} \leftarrow \text{iteration index}$$

3) fractional decay

$$\eta_t = \eta_0 / (1 + k t)$$

\* Additional methods:

- Newton's method

- compute learning rate instead of specifying it
- second order derivatives (Hessian matrix)

- AdaGrad:

· replaces Hessian computation (expensive) with approximations

- RMS Prop:

· add decay to ~~AdaGrad~~ AdaGrad

- ADAM:

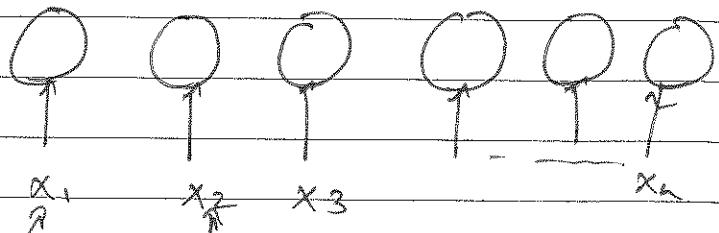
· add gradient momentum to RMS Prop

- Quasi-Newton methods (e.g. BFGS)

· approximate Hessian inverse (faster)

## Data ~~normalization~~

this gets more importance due to higher mag  $\downarrow$   $(100\dots1000)(0\dots)$



- features with high values have more influence and different features require a different learning rate.  $f_1, f_2, \dots, f_m$

- Data normalization

$$\bar{x}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\hat{\sigma}_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \bar{x}_j)^2$$

$$\tilde{x}_j^{(i)} = \frac{x_j^{(i)} - \bar{x}_j}{\hat{\sigma}_j}$$

normalized features

Ex 1	T
Ex 2	T
Ex 3	T
⋮	⋮
Ex m	T

avg  $[\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n]$

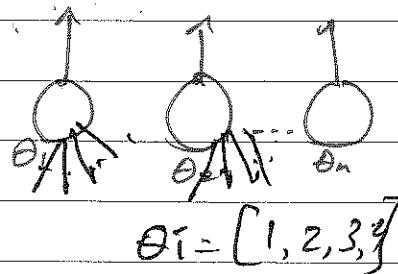
std  $[\hat{\sigma}_1, \hat{\sigma}_2, \dots, \hat{\sigma}_n]$

# Weight Initialization

$\theta = \text{constant} \rightarrow$  all outputs identical

$\rightarrow$  no learning

$\rightarrow$  random initialization

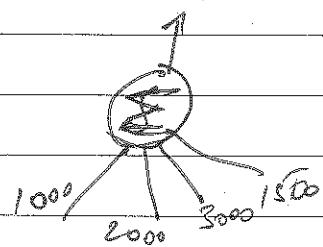


$\theta = \text{large} \rightarrow$  Saturated activation

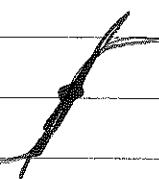
$\rightarrow$  Zero gradients

$\rightarrow$  no learning

$\rightarrow$  ~~initiated~~ initialize close to zero



$\rightarrow$  Initialize at random close to zero



data  
matrix

\* Glorot (Xavier) normal initialization

- Draw weights from normal distribution with:

$$\text{var}[w_i] = \left( \frac{\sigma^2}{\text{fan in} + \text{fan out}} \right)^{1/2} \quad \mu = 0$$

↓                      ↓  
Inputs                  Outputs

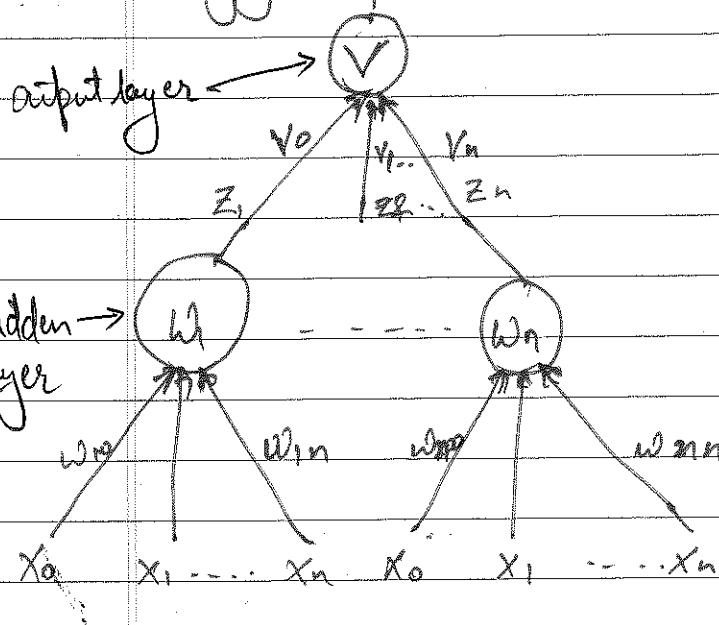
\* Glorot uniform initialization

- Draw weights from uniform distribution

within  $[-\text{lim}, \text{lim}]$      $\text{lim} = \sqrt{\frac{6}{(\text{fan in} + \text{fan out})}}$

# Back propagation :

$$E = \text{loss}(\hat{y}) \\ = \sum (y_j - \hat{y}_j)^2$$



Chain rule:

loss

$$\frac{\partial E}{\partial V} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial V} \rightarrow \text{find } V$$

Parameters (vector)

loss

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_j} \frac{\partial z_j}{\partial w_j} \rightarrow \text{find } w_j$$

Parameter (vector)

- Start with guess for parameters:  $V, \{w_j\}$  (at random or set to zero)

- Use  $x^{(i)}$  and current parameters to complete outputs:  $z^{(i)}, \hat{y}^{(i)}$

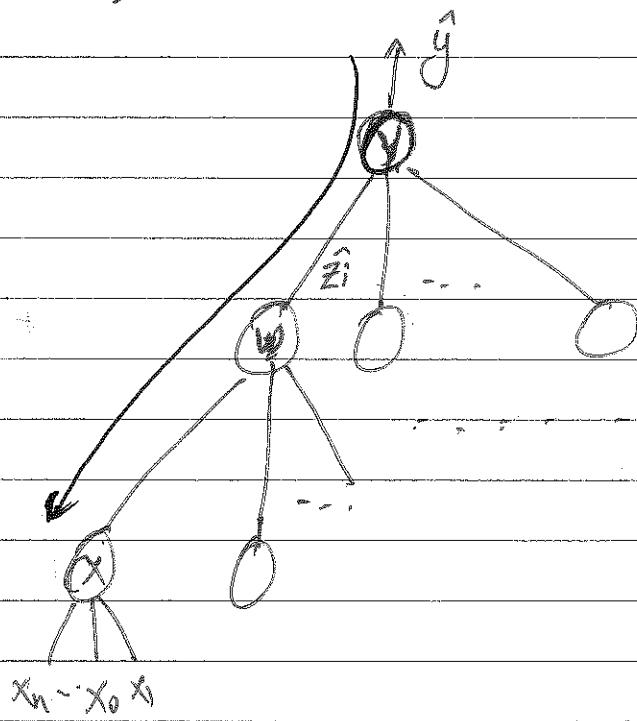
forward  
~~backward~~ pass  
push example

- Use outputs  $z^{(i)}, \hat{y}^{(i)}$  to update parameters:  $V, \{w_j\}$

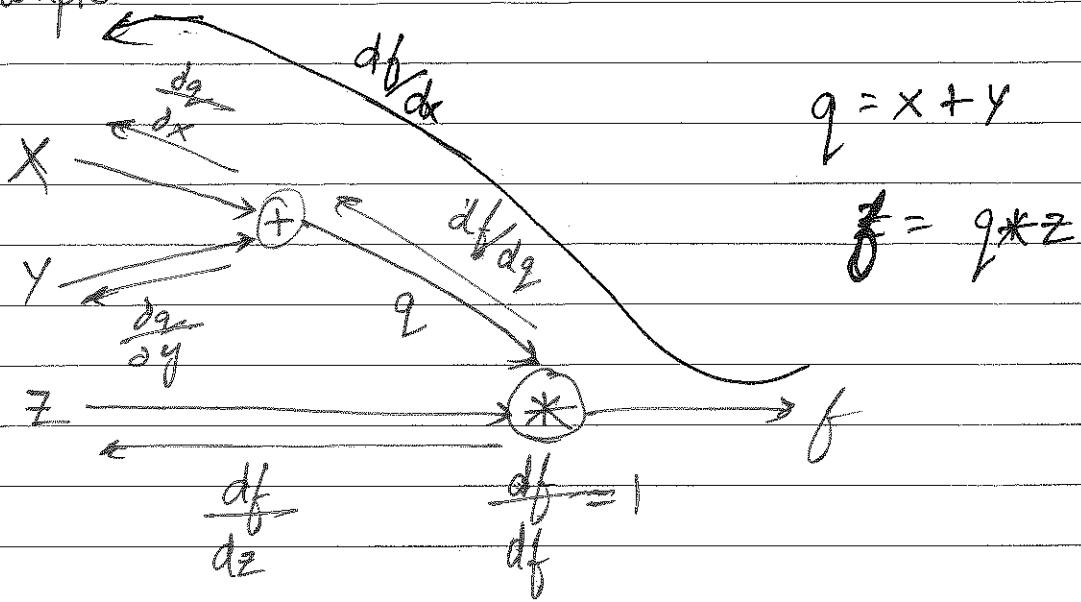
backward pass  
(push back gradients)

- continue while loss changes

- Represent network using a computational graph
- Because each node is simple, it has a simple explicit expression for its derivatives
- Forward pass: push input to compute all intermediate node values.
- Backward pass: starting with end nodes push gradients towards the beginning
- ~~Multiplying~~ Multiply backpropagated gradients (from back) by current gradients and propagate this



Example:



Want:  $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$



$$\frac{df}{dx} = 1$$

$$\frac{df}{dy} = z$$

$$\frac{dg}{dx} = 1$$

$$\frac{df}{dz} = g$$

$$\frac{dg}{dy} = 1$$

→ forward pass: push inputs to compute intermediate node values.

$$\text{eg. } (x, y, z) = (1, 2, 3)$$

$$g = 1 + 2 = 3$$

$$f = 3 \times 3 = 9$$

→ backward pass: push inputs to compute intermediate node values starting at end push gradients towards beginning.

propagate computed by node  
↓

$$\frac{df}{dx} = \frac{df}{df} \cdot \frac{df}{dg} \cdot \frac{dg}{dx} = 3 \cdot 1 = 3 \quad \frac{df}{dg} = z = 3$$

$$\frac{df}{dy} = \frac{df}{df} \cdot \frac{df}{dg} \cdot \frac{dg}{dy} = 3 \cdot 1 = 3$$

$$\frac{df}{dz} = g = 3$$

## Regularization Methods (prev. covered)

Large number of parameters tend to overfit.  
 ⇒ Need to prevent overfitting

### \* Dropout :

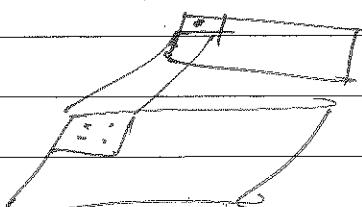
- at each training stage drop out units in fully connected layers with probability of  $(1-p)$ , where  $p$  is hyperparameter.
- Removed nodes are ~~are~~ reinstated with original weights ~~in~~ in the subsequent stages.

### \* Drop connect :

- Instead of dropping out unit, drop inputs with probability of  $1-p$  (i.e. randomness in inputs instead of output in normal dropout).

### \* stochastic pooling :

- A 'dropout' version for convolution layers.
- Activation within each pooling region is picked randomly according to a multinomial distribution.
- Equivalent to normal pooling from a set of randomly deformed images (deformation increases exponentially with layers).

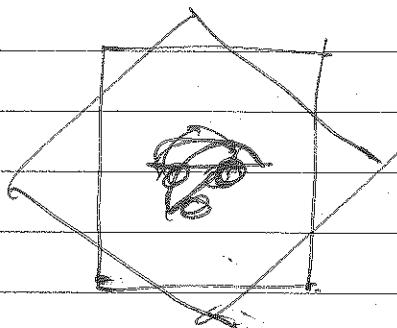


### \* Synthetic data : (data augmentation)

- increase variability in training data :

- synthesize new data

- Perturb existing data (eg crop images in different ways).



### \* Early stopping :

- stop training before learning completes.

### \* Reduce parameters :

eg: reduce hidden layers.

### \* Weight decay :

- multiply each coefficient by  $\rho \in [0, 1]$

As iterations progress weights that are not reinforced decay to 0.

$$(\rho)^{100} \rightarrow 0$$

- Equivalent to adding regularization term to loss function:

\* Max norm constraints:

- clamp to weights vector at each unit so that  $\|w\|_2 \leq c$

\* Batch normalization