

# SMA I (07/01/2020)

- collect some data - D.  $p$ -dim. vector  $x_1, x_2, \dots$  → label space  $\{1, \dots, 10\}$ .
- input  $p$ -dim. vector and output is a value - classification problem.
- ML starts with collecting good data [Imp].
- Goal - to do well on collected data?  
OR, to do well on future data?
- Goal is not to do perfect on collected data but to do good on future data.  
Overfitting on current data is not good.



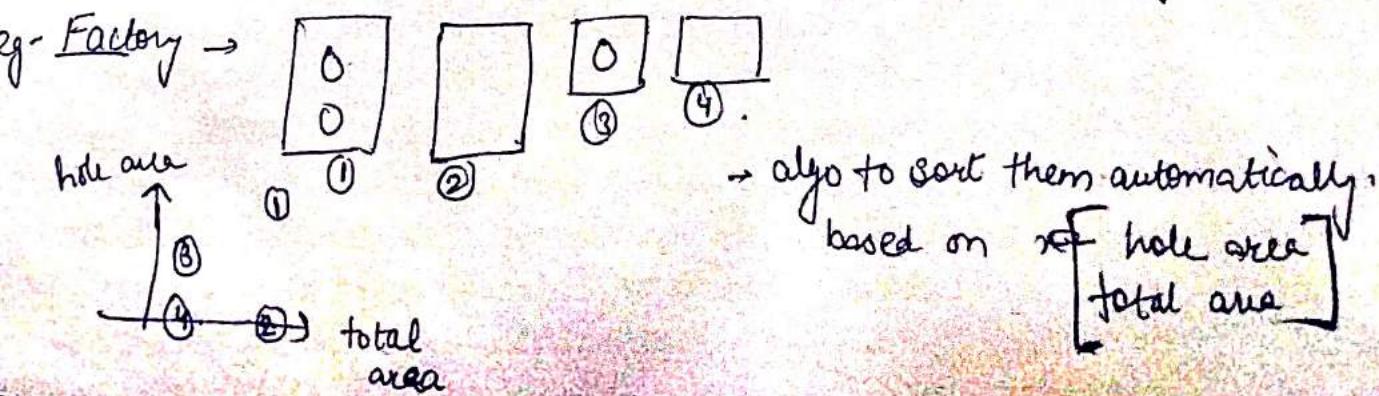
- Overfitted one - has less loss percentage  
[loss percent = diff. b/w predicted & actual value].

- Data — 

Train	Test
-------	------

 making data splits →  
[very careful during splitting data].

- K-nearest Neighbours:  
 $k$  - odd no. (generally).  
major drawback - we have to calculate dist. from all points in dataset & then find  $k$ -nearest for classification.  
This ~~causes~~ dist. calculation is very time consuming.



$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad x' = \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix}$$

$$\text{euclidian distance} = \sqrt{(x_1 - x'_1)^2 + (x_2 - x'_2)^2}$$

$\Rightarrow$  Cosine distance - Angle b/w two objects

# document - collection of words

$\rightarrow$  bag of words representation.

$[x \text{ has a macho figure. In laboratory test.}]$

doc 1

$[$  figure 1 illus. our study laboratory test forms  $\dots$  We assume labo  $]$

doc 2

Suppose in our vocabulary, there are 2 words - laboratory, figure.

we represent doc. by bag of words -  $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix} \rightarrow$  sparse generally.

↑  
bag of Many values may be 0.  
words vector.

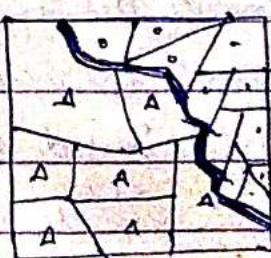
$\rightarrow$  According to data, we have to decide which algo to use, which distance to use (euclidian or cosine).

### Nearest Neighbour classification

$\rightarrow$  Voronoi Tessellation

$\rightarrow$  partition space in regions.

$\rightarrow$  boundary : points at same dist from two diff. training examples.



[1-N classification]

$\rightarrow$  We look at one point near the data and based on it decide its region.

→ decision boundary : non-linear ; compare with other classification approaches.

not a single line - non-linear, arbitrary.

→ Issue : sensitive to outliers.

Single mislabeled example can significantly change boundary.

### KNN Classification

→ Given pairs of training examples  $\{x_i, y_i\}$ , we want to classify a testing point  $x$ .

→ Algo -

(i) Compute distance  $D(x_i, y_i)$   $D(x_i, x_j)$  to every training example  $x_j$ .

(ii) Select  $k$  closest instances  $x_1, \dots, x_k$  and their labels  $y_1, \dots, y_k$ .

(iii) Output the class  $y^*$  which is most frequent in  $y_1, \dots, y_k$ .

→ No parametric approach, no training.

### KNN fn. learning

→ Given pair of training examples  $\{x_i, y_i\}$ , we want  $y_i \in R$  for given input  $x$ .

→ Algo - (i) Compute distance  $D(x_i, x_j)$  to every training example  $x_j$ .

(ii) Select  $k$  closest instances  $x_1, \dots, x_k$  and their labels  $y_1, \dots, y_k$ .

(iii) Output the weighted mean of  $y_1, \dots, y_k$ .

weight  $\propto 1/\text{distance} \rightarrow$  closer : more weight & farther : less weight.

[This algo is mainly used for classification but can be used for regression as well.]

How to decide value of  $k$ ?

- if value of  $k$  is large — the majority class will always be the result i.e. class having minority won't be selected in result.
- if  $k$  is very small — due to some noise, the result highly susceptible to outliers.
- need a decent  $k$ .
- experiment with diff.  $k$ ; dividing dataset; cross-validation.

→ divide dataset into 3 parts

train	10%	Test
80%	Validation	10%

train with 80% data and the test for  $k=1-n$  for the validation data; decide on what  $k$  are you getting max. accuracy (how many 10% can you correctly classify).

Test on testing data.

⇒  $k$  should always be less than  $\sqrt{n}$  {  $k \leq \sqrt{n}$  } .  
{ 5, 7, 11 commonly used? } usually take much lower than that.

→ checking  $k$  on test data — overfitting

→ Distance Metric.

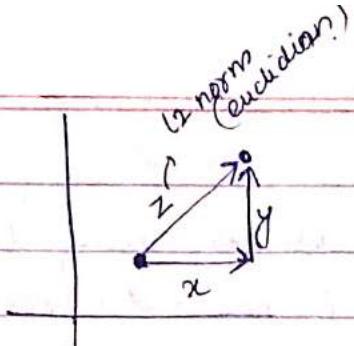
$$(a) \text{ Euclidean } D(x, x') = \sqrt{\sum_i (x_i - x'_i)^2}$$

→ (dividing) splitting dataset : split happens across labels also (each label 80%); signature verification; handwriting — all data must not be in training. and some should be left for test untouched.

→ Mahalanobis distance :  $D(x, x') = (x - x')^T \Sigma^{-1} (x - x')$ .

→ Minkowski Distance Norm ( $p$ -norm) :  $D(x, x') = \sqrt[p]{\sum_i (x_i - x'_i)^p}$

- $p=1 \rightarrow$  manhattan
- $p=2 \rightarrow$  euclidean dist.
- $p=\infty \rightarrow$  intuiting read about them.



$$x + y = z \\ \downarrow \\ \text{mahalanobis dist} - \\ \boxed{\text{L2 Norm}}$$

→ cosine Dist.

→ Hamming Dist.

$$D(x, x') = \sum_i 1_{x_i \neq x'_i}$$

eg -  $\begin{cases} @at & D=1+0+0. \\ @at & \end{cases}$

→ KL Divergence.

$$D(x, x') = \sum x_i \log \left( \frac{x_i}{x'_i} \right)$$

→ information, entropy, mutual info, self energy.

→ comes from information theory.

→ not symmetric.

Norm of a vector [Not 2 vectors]

Norm of a vector [Not 2 vectors]

## SMAI (lecture - 2) (11/01)

Minkowski Dist. → if  $p = \infty$ , it selects  $\max_i (x_i - x'_i)$ .

$$\left( (x_0 - x'_0)^p + (x_1 - x'_1)^p + (x_2 - x'_2)^p + \dots \right)^{1/p}$$

$$\text{eg} - \left( (3)^{1000} + (1.5)^{1000} + (4)^{1000} \right)^{1/1000} \approx (4)^{1000} \approx 4.$$

→ in comp. to other values max. value will dominate in case of  $p = \infty$ .

$$\# \text{ when } p = \infty \rightarrow \boxed{\max_i (x_i - x'_i)}$$

# When  $p=0 \Rightarrow$

→ Maddam Ellipses - is a region on a chromaticity diagram which contains all colors which are indistinguishable to avg. human eye from the color at the centre of ellipse.

→ How to pick odd  $K \rightarrow$  flip coin, pick class with greater prior

→ Suppose we do 7-NN. and there are 4 classes  $c_1 c_2 c_3 c_4$ .  
 $\begin{matrix} & c_1 & c_2 & c_3 & c_4 \\ 3 & & 3 & 1 & 0 \end{matrix}$ .

Now to choose from  $c_1$  &  $c_2$  - we use 1-NN. i.e. what is closest is  $c_1$  &  $c_2$ .

→ what if there are missing values?

take median, mean etc. - replace by some value. performing calc on that class.

Applications of KNN → manual data entered to be read )  
these all are images.

① MNIST dataset

→ decent performance when lots of data

→ Nearest neighbour is competitive.

5	0	4	1
9	2	1	3

[handwriting  
recognition]

② Problem : Where (which country or GPS location) was this image

↑ taken?

Research Paper.

# Data is a challenge ; representation is also a challenge.

→ image - set of pixels - 2D grid.

how problem 2 was solved.

- (i) Get 6M images from Flickr with gps info (dense sampling across world)
- (ii) Represent each image with meaningful features.
- (iii) Do kNN (large  $k$  better, they use  $K=120$ )

③ Given an image can you caption it?

→ query expansion approach  $q = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M \text{sim}(I_q, I_j) \cdot c_i$

# Read word2vec.

## KNN Pros & Cons

- Simple, easy to understand.
- Non parametric (No training req.)
- Widely Applicable.
- No assumption about data.
- Easy to update.
- Could be used as baseline algo.
- K & D - only two choices.

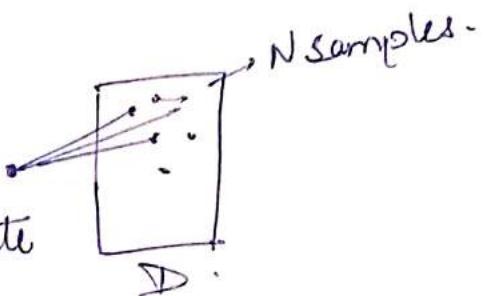
### Cons

- computation time is very high. - the major issue.
- sensitive to outliers.

→  $O(n \cdot d)$  — Time Complexity.

d - dimensionality

for each of the sample, you compute N times.

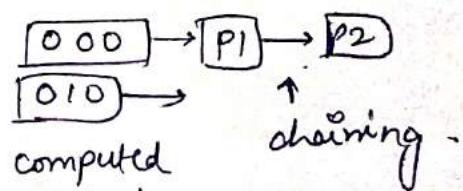
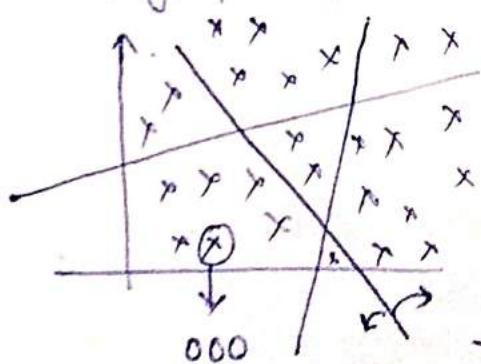


## Making KNN Fast

- Ideas :
- (i) Reduce d - dimensionality reduction.
  - (ii) Reduce n - don't compare with all training example.
    - Kd trees  $O(d \log_2 n)$
    - locally sensitive hashing.
    - inverted indices.

### Locally Sensitive Hashing (- Really Useful)

- Line / Hyperplane (4-D & above) divide Space in 2 parts.



- define k-hyperplanes, randomly decided.

Given  
→ for a test point we compute hash and then apply knn in the chain of that hash only.

for 1 plane, 2 choices, k-planes -  $2^k$  choices.  
distance from  $n/2^k$  points.

- point → k-dim. vector.

→ Complexity →  $O(kd + dn/2^k)$

k-dim, d-planes & then checking +ve / -ve for ~~all hyperplanes~~ hashing.

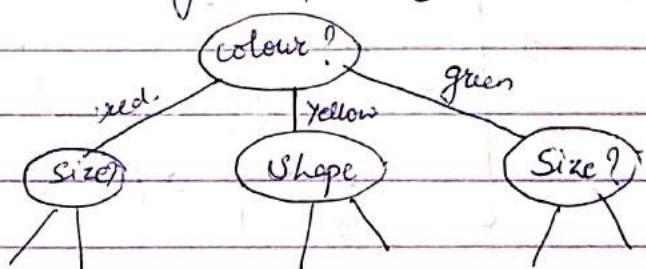
⇒ Parse Window — anything within a const. dist. will be considered.  
is nearest neighbours.

SMAI (18/01/2020) .

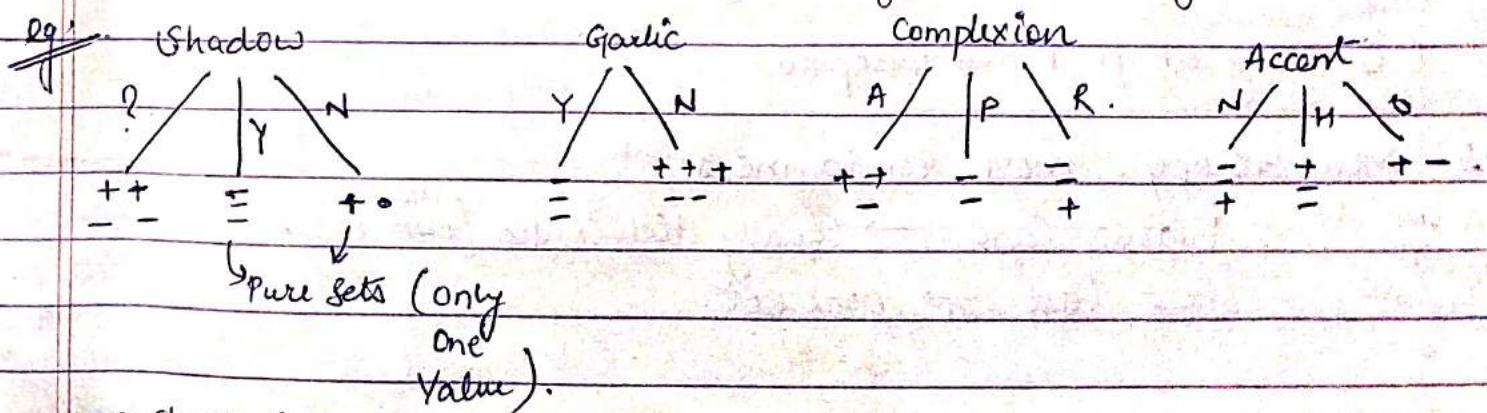
### Decision Trees

- one of the most intuitive classifiers.
- easy to understand and construct.
- Surprisingly, also works very (very) well.
- Netflix used decision tree for recommendation feature. (had large data but algo worked well).

Eg:- fruit described using 4 tuple : {color, shape, taste, size}.



# Goal : to build this tree ; data (large amount) is given.



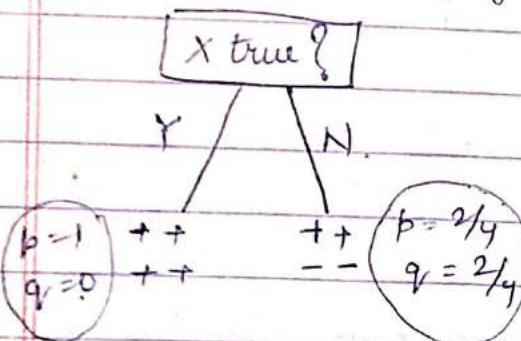
→ choose factor having more pure sets.

Hence Shadow → Garlic (for decision tree)

## Impurity functions

- ① Entropy :  $i(v) = -[q \log q + (1-q) \log(1-q)]$  ] More commonly used
- ② Gini Index :  $i(v) = 2q(1-q)$
- ③ Misclassification Rate :  $i(v) = \min(q, 1-q)$

$$\begin{aligned} \text{Entropy} \Rightarrow & \sum -p_i \log p_i \\ & = -(p \log p + (1-p) \log(1-p)) \rightarrow \text{Since 2 classes.} \end{aligned}$$

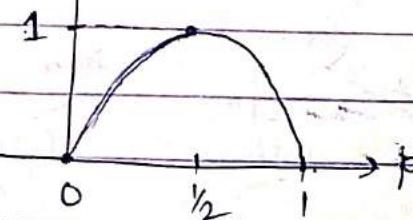


Set 2 entropy.

$$\begin{aligned} & = -\left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2}\right) \\ & = 1 \end{aligned}$$

$p = +ve$  class

$q = -ve$  class.



entropy set 1.

$$= -(\log 1 + 0 \log 0)$$

$$= 0.$$

$$\left\{ \begin{array}{l} \lim_{p \rightarrow 0} p \log p = \lim_{p \rightarrow 0} \frac{p}{-1/p^2} = 0 \end{array} \right\}.$$

Similarly at  $p=1 \Rightarrow \text{entrop}=0$ .

# More Entropy, more randomness

Neither sure which way, both equally likely  $\rightarrow$  Worst Case  $\rightarrow$  Equally distributed [Most random].  
Best Case  $\rightarrow$  pure set.

# Most random thing will be equally distributed

# Sets which gives ~~higher~~ lower entropy are better.

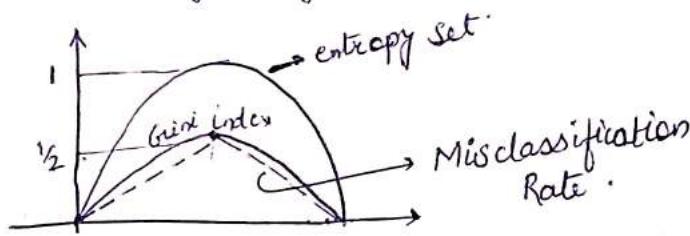
Shadow  
 $\frac{1}{2}(\text{entropy } S_1) + \frac{3}{8}(\text{entropy } S_2) + \frac{1}{8}(\text{entropy } S_3)$   
 $= \frac{1}{2}(1) + \frac{3}{8}(0) + \frac{1}{8}(0) = 0.5$

$S_1 \quad | \quad S_2 \quad S_3$   
 $\underline{\underline{++}} \quad \underline{\underline{=}} \quad +.$

Entropy [Gini] =  $\frac{3}{8}x_0 + \frac{5}{8}x_-$  = (0.6) Suppose.

Entropy [Complexion] =  $\frac{3}{8}x_- + \frac{2}{8}x_0 + \frac{3}{8}x_-$

- ⇒ Entropy Set = weighted sum of entropies of subsets
- ⇒ In real world, entropy set & gini index are used.
- ⇒ Check Entropy diagrams. (Slides).



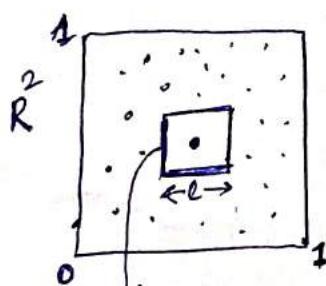
- ⇒ In case of numeric data — we distinguish on the basis of numerical value.  
[ $x > x_1$ ;  $x < x_2$  — etc.]

# When do we Stop?

We use validation sets [Best Way].  
→ more partition will lead to overfitting.

SMAI (21/01/2020)

### Curse of Dimensionality



that should be the size on avg., I can have  $k$  points.

[Total area = 1].

let,  $N$  points ( $\leq 1000$ ) in a sq. region —

find  $k$  nearest neighbours ( $k=10$ ).

→ uniformly distributed (not dense or sparse at any pt.).

$$\frac{k}{N} \approx l^2$$

$$l \approx \left(\frac{k}{N}\right)^{1/2} \rightarrow 2 \text{ dimensions.}$$

$$\frac{k}{N} \approx l^d \rightarrow d \text{ dimensions.}$$

hypercube

$$l = \left(\frac{k}{N}\right)^{1/d}.$$

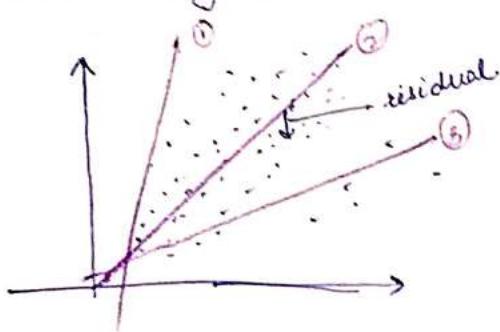
KNN assumption — points which are closer are more likely to be part of same class.

$d$	$e$	$\rightarrow N = \frac{k}{d} \text{ if } l = 10 \Rightarrow [N = k \cdot 10^d]. \text{ data points req.}$
2	0.1	for 10% of space.
10	0.63	
100	0.995	
1000	0.9954	

→ As  $d$  increases, area req. to find  $k$  points increases.  
 → More amount of data req. for training.  
 Hence, called "curse" of dimensionality.

- ⇒ But, still KNN is used & give good result as data is not uniformly distributed. Data lies on many fold.
- ⇒ Sometimes, KNN won't work, this may be the reason.

## Linear Regression



which line describes the data best?

- find sum of distance of points from line -  
 The line from which sum is least is best line.
- find a line, sum of errors is minimum.

$$y = w^T x.$$

$$y = w_0 + w_1 x_1 + w_2 x_2^2 + w_3 x_3^3 - \dots \text{-- w ard.}$$

Still,  $y = w^T x.$

- Output can be predicted as weighted sum of ~~distances~~ inputs.

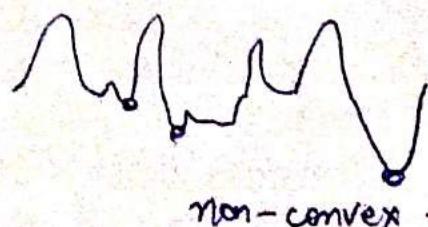
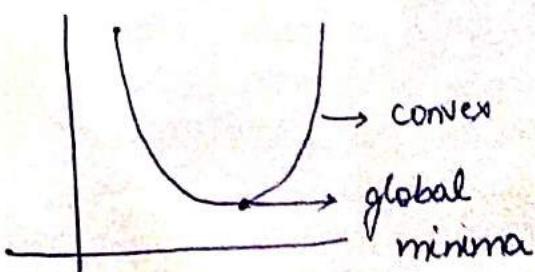
$$y = mx + c$$

$w \rightarrow$  set of parameters  
vector.

$$J(w) = \sum_i e_i^2 = \sum_i (w^T x_i - y_i)^2 \rightarrow \text{cost function}$$

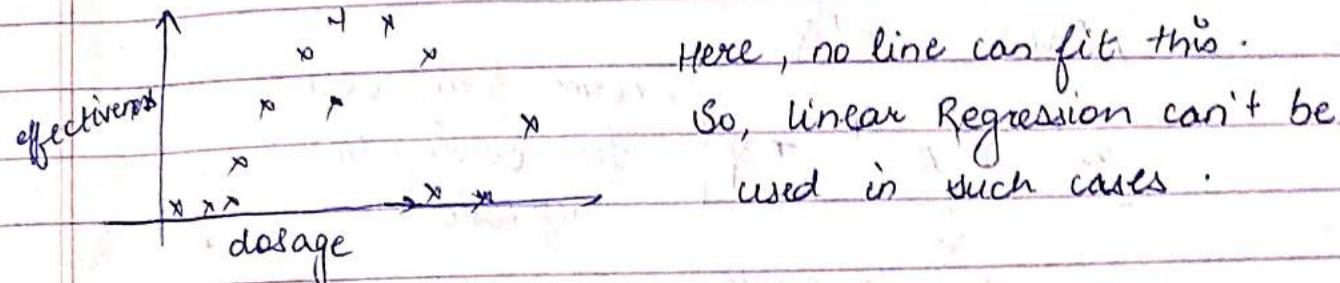
- our task is to minimize this cost fn.

- Eventually, we have to find set of  $w$ .



- polynomial fitting can be done using linear regression.

- $x$  needs to be linear,  $y$  can be anything.  
 for  $x_i \rightarrow$  output  $y_i$ :  
 $\sum_{i=1}^n (w^T x_i - y_i)^2 \rightarrow$  error to be minimized.
- ⇒ Where Linear Regression can't be used?



→ Solve for  $J(w)$  which minimizes  $J(w)$ .

$$J(w) = \frac{1}{2} \sum_i (w^T x_i - y_i)^2$$

$$A = \begin{bmatrix} -x_1^T \\ -x_2^T \\ -x_3^T \\ \vdots \\ -x_n^T \end{bmatrix}_{N \times d}$$

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{d-1} \end{bmatrix}_{d \times 1}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$

$$J(w) = \frac{1}{2} (AW - Y)^T (AW - Y)$$

$$AW = \begin{bmatrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_n^T w \end{bmatrix}$$

$$\begin{bmatrix} 1 & x_1^T & x_2^T & \cdots & x_d^T \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{d-1} \end{bmatrix} = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_d x_d = w^T x_i.$$

$$[\text{Hence, } AW = w^T x_i]$$

$$\{(AW - Y)^T (AW - Y)\} = (AW - Y)^2$$

$$\frac{\partial J}{\partial W} = \frac{1}{2} \frac{\partial}{\partial W} [(AW - Y)^T (AW - Y)]$$

$$= \frac{1}{2} \frac{\partial}{\partial W} [(AW)^T AW - (AW)^T Y - Y^T AW + Y^T Y]$$

$$P = AW ; Q = Y , P^T Q = Q^T P$$

$$= \frac{1}{2} \frac{\partial}{\partial w} \left[ (Aw)^T Aw - 2(Aw)^T Y + Y^T Y \right].$$

$$= \frac{1}{2} \frac{\partial}{\partial w} \left[ w^T A^T Aw - 2w^T A^T Y + Y^T Y \right].$$

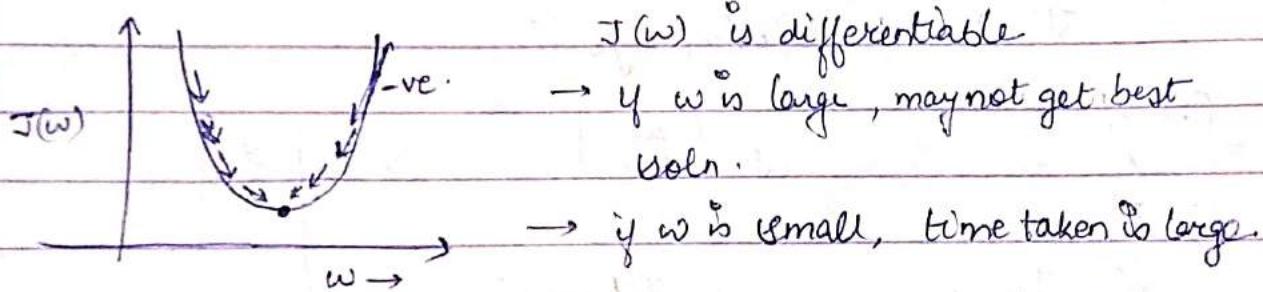
$$= \frac{1}{2} \frac{\partial}{\partial w} \left[ 2A^T Aw - 2A^T Y \right]$$

$$\Rightarrow A^T Aw = A^T Y$$

$w = (A^T A)^{-1} A^T Y$

normal form soln.  
for solving linear  
regression.  
pseudo inverse.

Gradient descent



$\rightarrow$  we find gradient and go in negative direction to reach 0.

$$w_{t+1} \rightarrow w_t - \lambda J'(w_t).$$

Taylor Series :

$$\begin{aligned} J(w+h) &\approx J(w) + h J'(w) + \dots \\ &\approx J(w) + h J'(w) \end{aligned}$$

$$w_{t+1} = w_t - \lambda J'(w).$$

$$h = -\lambda J'(w)$$

$$J(w_t - \lambda J'(w)) \approx J(w) - (+\lambda J'(w)) J'(w).$$

$$J(w_{t+1}) \approx J(w) - \lambda (J'(w))^2.$$

$$|J(w_{t+1}) \leq J(w)|$$

# Gradient descent always work till curve is convex.

## SMAI (25/01/2020)

$$w_{t+1} \leftarrow w_t - \lambda J(w_t)$$

$\Rightarrow$  Amount of step depends on  $\lambda$  and  $J(w)$ .

$$J(w) = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$$

$$J'(w) = \sum_{i=1}^N x_i (w^T x_i - y_i)$$

Algo :  
① Choose a learning rate,  $\lambda \in R^+$

② Choose a random  $w$ ,  $w \in R^{d+1}$

③ Repeat ④ till converges.

$$④ w = w - \lambda \sum_{i=1}^N x_i [w^T x_i - y_i]$$

$\begin{cases} \lambda \text{ bigger} - \text{ bigger rate} \\ \lambda \text{ smaller} - \text{ smaller rate} \end{cases}$

$\leftarrow$  Batch gradient descent Algo.

$\rightarrow$  Problem - we have to go through all data once to make one step.

Repeat for all samples -

$$w := w - \lambda x_i (w^T x_i - y_i)$$

$\leftarrow$  stochastic  
gradient  
descent

$\hookrightarrow$  Problem :  $w$  will fluctuate a lot ; data is noisy. So cost fluctuates a lot.

### Mini Batch Gradient Descent

Batches of samples from the training set (eg - 32 samples)

$$w = w - \lambda \sum_{i=1}^{32} x_i (w^T x_i - y_i)$$

finish one pass over the entire data } epoch.

\* Extrapolation - don't always work.

\* Overfitting - avoid it.

### Regularization (Imp.)

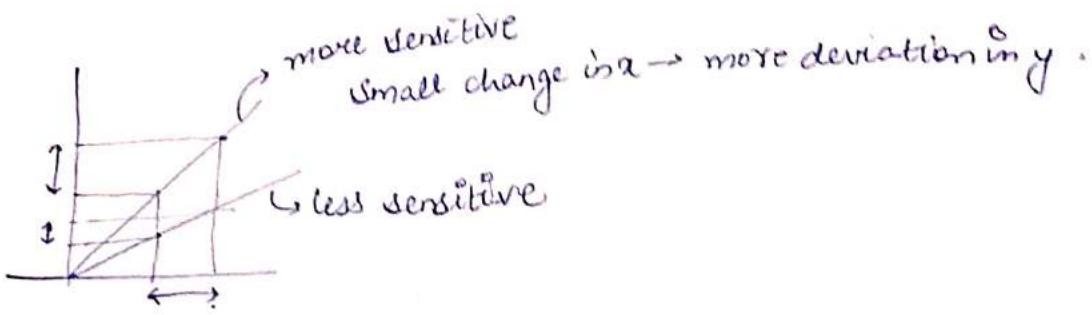
$$J(w) = \sum_{i=1}^N (w^T x_i - y_i)^2 + \left( \sum_{i=1}^d w_i^2 \right)$$

Norm of weights  
 $\|w\|_p = (\sum |w_i|^p)^{1/p}$

$[J'w$  can still be computed  
efficiently.]

$$J(w) = \sum_{i=1}^N (w^T x_i - y_i)^2 + \lambda \sum_{i=1}^d \|w_i\|_2$$

$\rightarrow$  for  $L_2$  norm.



$$J(w) = \sum_{i=1}^N (w^T x_i - y_i)^2 + \lambda \sum_{i=1}^d \|w_i\|_2.$$

$$\boxed{y_i = w_0 + w_1 x + w_2 x^2 + \dots + w_{25} x^{25}}$$

very large than

without regularization

$\sum (w^T x_i - y_i)^2$	will be small
but other term is very large.	

⇒ With regularization, term 2 can't be large, so suppose it's modest, voln. is slightly worse but it's still small.

There is tradeoff.

⇒ Regularization helps reduce overfitting.

⇒ Ridge Regression = Normal Regression +  $L_2$  norm regularization.

⇒ Lasso Regression = Normal regression +  $L_1$  norm regularization.

↓

works as

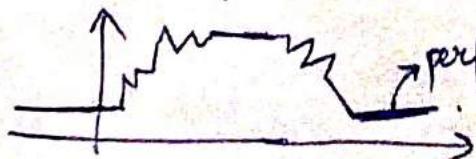
$$\boxed{J(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \sum_{i=1}^d |w_i|}$$

feature selection.

hyperparameters — found using cross-validation

{contour?}.

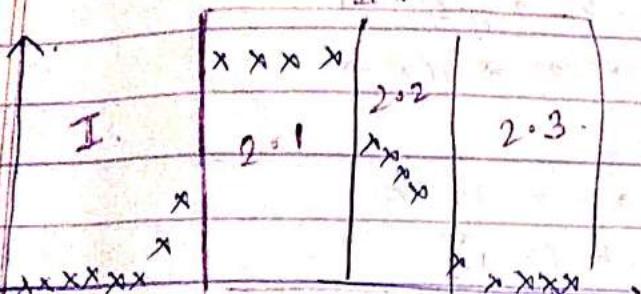
⇒ Lasso can give perfect slope but ridge doesn't give.



perfect slope \* ridge can't give this.

## Elastic Net Regression

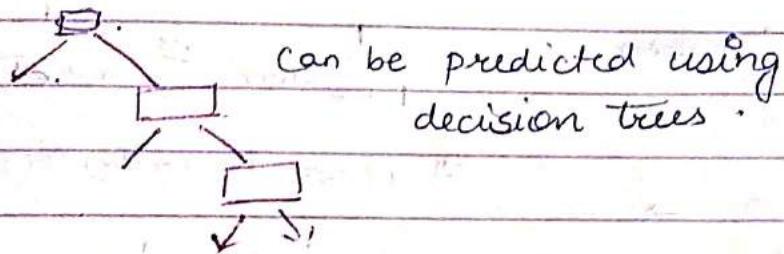
$$\sum (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda_1 \sum \|w_i\|_2 + \lambda_2 \sum \|w_i\|_1$$



Here regression can't be used. No curve can fit effectively.

we make partitions.

→ We use decision trees



SMAI (28/01).

## Logistic Regression

- classification approach

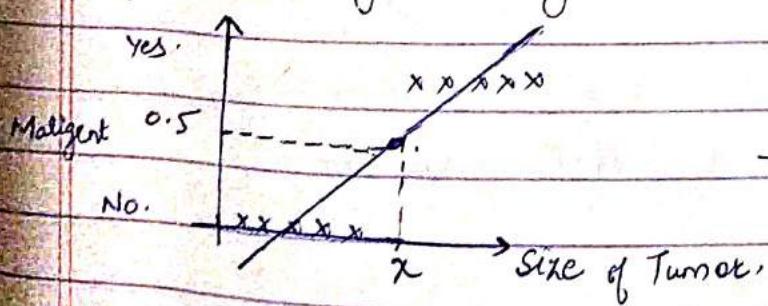
classification problems - (i) spam or not (2 class prob.)

(ii) Spam, primary, promotion (multi-class)

(iii) Tumor - malign or benign

(iv) Transaction - genuine or fraud.

logistic = logit = Sigmoid

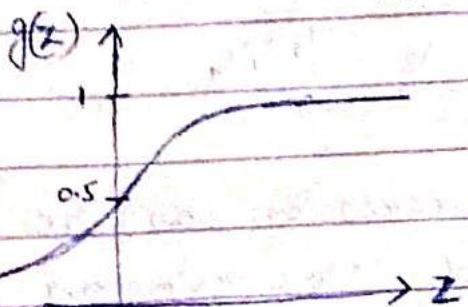


How to make this problem a classification problem?

→ we decide a threshold  $x$  and calculate  $y$  at it. on line.

Now for any  $x$ , if  $y(x) > \text{threshold}$ , we classify it as 'Yes' else classify it to 'No' label.

$$g(z) = \frac{1}{1+e^{-z}} \quad \leftarrow \text{Sigmoid fn.}$$



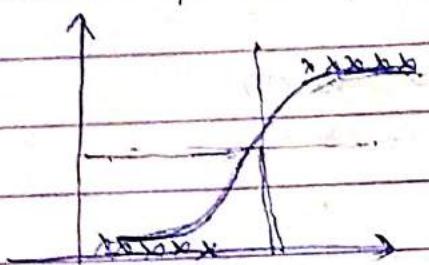
$$z=0 \Rightarrow g(z)=0.5$$

$$z \rightarrow \infty \Rightarrow g(z) \approx 1$$

$$z \rightarrow -\infty \Rightarrow g(z) \rightarrow 0.$$

Fn. nice. [ TR  $\rightarrow -\infty$  to  $\infty$ .  
output  $(0, 1)$ .

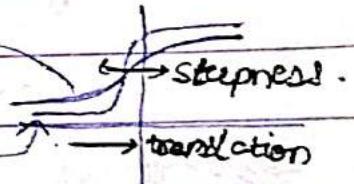
For our problem, we can make some shifted sigmoid fn.



$$\frac{1}{1+e^{-wx}} \approx \frac{1}{1+e^{(-\text{bias})}}$$

$w$  can be controlled

$\rightarrow$  if  $w$  is small — curve smooth  
 $w$  is large — curve steep



let,

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

$$\text{Then, } \sigma'(x) = \frac{+e^{-x}}{(1+e^{-x})^2} = \left(1 - \frac{1}{(1+e^{-x})}\right) \frac{1}{(1+e^{-x})} \\ = \sigma(x) \cdot [1 - \sigma(x)].$$

$$[\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))]$$

$\Rightarrow$  In sigmoid fn., we want to control stepness and translation mainly.

In linear regression  $\rightarrow J(w) = \sum_i (y_i - w^T x_i)^2$   
 logistic regression  $\rightarrow J(w) = \sum_i [y_i - \sigma(w^T x_i)]^2 \rightarrow$  not a convex function.

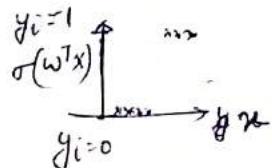
### Maximum likelihood Estimation

$$P(y_i = 1 | x_i, w) = \sigma(w^T x_i)$$

$$P(y_i = 0 | x_i, w) = 1 - \sigma(w^T x_i)$$

2 class classification.  
 $\left\{ \begin{array}{l} y_i \in 0 \text{ or } 1 \text{ as it can belong} \\ \text{to a class or not, no is blur} \end{array} \right\}$

$$\frac{P(y_i | x_i, w)}{\text{likelihood}} = \underbrace{(\sigma(w^T x_i))^{y_i}}_{y=0 \times} \times \underbrace{(1 - \sigma(w^T x_i))^{(1-y_i)}}_{y=1 \times}$$



Maximize this - our goal

bcz. when  $y_i = 0$ , we want  $\sigma(w^T x_i)$  to be close to 0.

$\rightarrow$  To maximize  $\sigma(w^T x_i)$  term  $\rightarrow$  minimize  $(1 - \sigma(w^T x_i))$ . as at  $y_i = 0$ , 1st term doesn't come into play  
 at  $y = 1 \Rightarrow$

# Maximize P  $\Rightarrow$  Is it a right goal?

$y_i = 0 \rightarrow \sigma(w^T x)$  is close to zero

$y_i = 1 \rightarrow \sigma(w^T x)$  is close to 1

So, by maximizing single fn. P  $\rightarrow$  we achieve our goal. Hence, a good idea.

$$L(w) = \prod_{i=1}^N P(y_i | x_i, w) \rightarrow \text{maximize this.}$$

$\rightarrow$  finding parameters for w  $\rightarrow$  turns to  $\rightarrow$  maximizing above fn.  
 $\rightarrow$  find set of w where above fn. is maximized, we are done.

$$\log(L(w)) = \log \prod_{i=1}^N P(y_i | x_i, w)$$

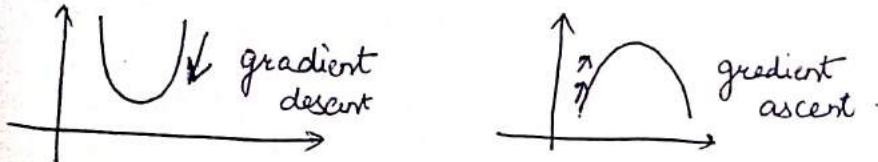
$$G(w) = \log(L(w)) = \sum_{i=1}^N \log(P(y_i | x_i, w)).$$

$$= \sum_{i=1}^N \log((\sigma(w^T x_i))^{y_i} \times (1 - \sigma(w^T x_i))^{(1-y_i)})$$

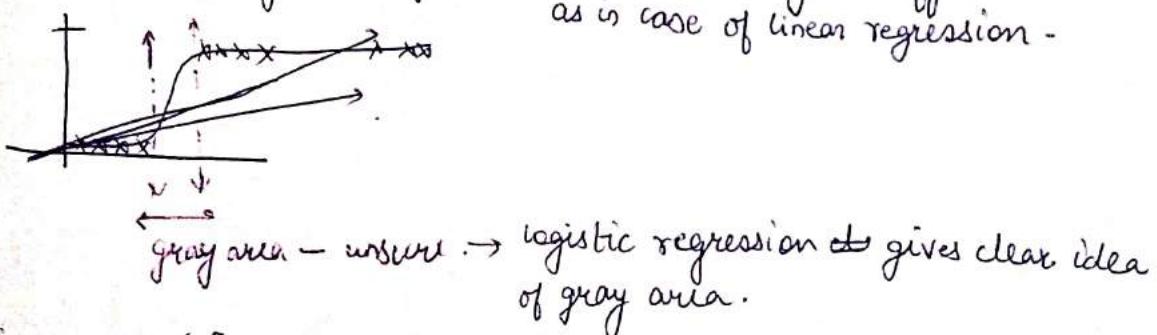
$$G(w) = \sum_{i=1}^N y_i \log(\sigma(w^T x_i)) + (1-y_i) \log(1 - \sigma(w^T x_i)).$$

$$\begin{aligned}
 G'(\omega) &= \sum_{i=1}^N y_i \cdot \frac{1}{\sigma(\omega^T x_i)} \cdot \sigma'(\omega^T x_i) \cdot x_i + \frac{(1-y_i)}{1-\sigma(\omega^T x_i)} \sigma'(\omega^T x_i) x_i \\
 &= \sum_{i=1}^N \left( \frac{y_i}{\sigma(\omega^T x_i)} - \frac{(1-y_i)}{1-\sigma(\omega^T x_i)} \right) \times \sigma'(\omega^T x_i) \cdot x_i \\
 &= \sum_{i=1}^N \left( \frac{y_i}{\sigma(\omega^T x_i)} - \frac{(1-y_i)}{1-\sigma(\omega^T x_i)} \right) \sigma(\omega^T x_i) (1-\sigma(\omega^T x_i)) \cdot x_i \\
 &= \sum_{i=1}^N (y_i(1-\sigma(\omega^T x_i)) - (1-y_i)\sigma(\omega^T x_i)) x_i \\
 &= \sum_{i=1}^N (y_i - y_i \cdot \sigma(\omega^T x_i) - \sigma(\omega^T x_i) + y_i \sigma(\omega^T x_i)) x_i \\
 G'(\omega) &= \sum_{i=1}^N (y_i - \sigma(\omega^T x_i)) x_i
 \end{aligned}$$

→ Goal was to ~~maximize~~ maximize  $G(\omega)$ .



→ In case of logistic regression, outliers may not affect much.  
as in case of linear regression.



{ After training,  $\sigma(\omega^T x) > 0.5$  - class 1  
 $\sigma(\omega^T x) \leq 0.5$  - class 2. }

### 2 Class classification

one vs one

4 classes  $\rightarrow 4C_2$  classifiers.

A, B, C, D

A vs B

A vs C

A vs D

B vs C  
B vs D  
C vs D

### Multi class classification

one vs all.

$C_1 \rightarrow A$  vs B, C, D.

$C_2 \rightarrow B$  vs A, C, D.

$C_3 \rightarrow C$  vs A, B, D

$C_4 \rightarrow D$  vs A, B, C.

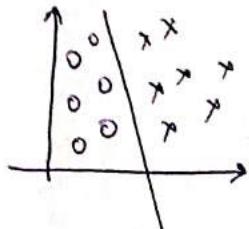
run all &

take decision

SMA1 (01/02/2020)

## SVM (Support Vector Machines)

→ linearly separable data (assumption)

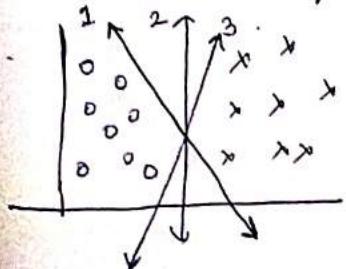


linear classification in 2-D : learning line (2 parameters)  
in 3-D : plane (3 parameters).

(We assume no noise for now)

→ Multiple solns. exist for linearly separable data.

Are all solns. equally good?

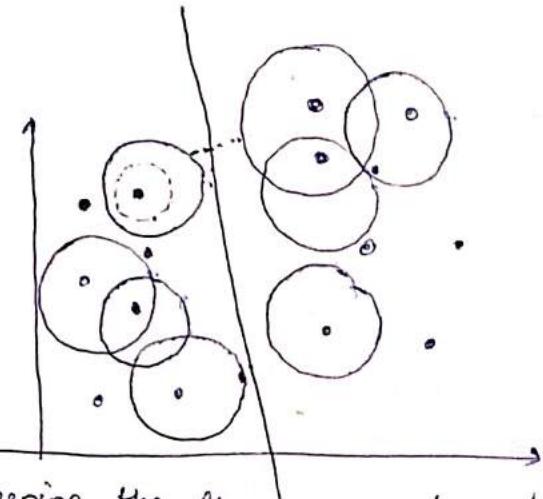


→ line 2 is better as it gives more margin.

Margin : bubbles around samples.

Margin - radius of a region around each training sample, through which the decision boundary can't pass.

As margin increases, the feasible region reduces.



- Keeping the line away from bubble
- increasing size of bubble.

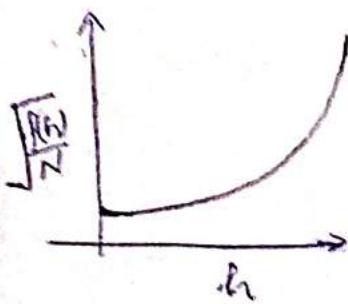
Goal : line which divides data in two parts (assuming linearly separable data). We have to find a line which divides data and maximizes margin. [Algo called SVM].

Bound on Expected loss

$$= R(x) \leq R_{\text{train}}(x) + \sqrt{\frac{f(h)}{N}}$$

$h$  is VC dimension and  $\{ f(h) = h + h \log(2N) - h \log(h) - c \}$ .

- monotonically increasing fn.  
 → taking lower value of  $h$  gives lower value of  $\sqrt{f(h)/N}$

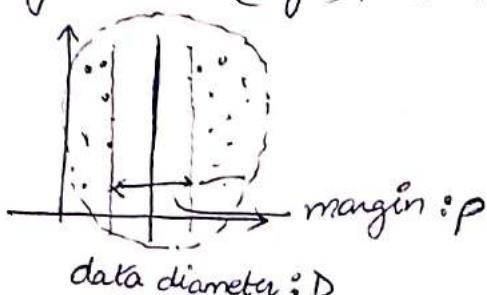


### Why maximise the margin?

- To reduce test error, keep training error low (say 0), and minimize the VC dimension,  $h$ .

$$\text{Relative margin} = \frac{P}{D}$$

$$\text{VC} = D, h \leq \min\left\{\alpha, \left\lceil \frac{D^2}{P} \right\rceil\right\} + 1$$

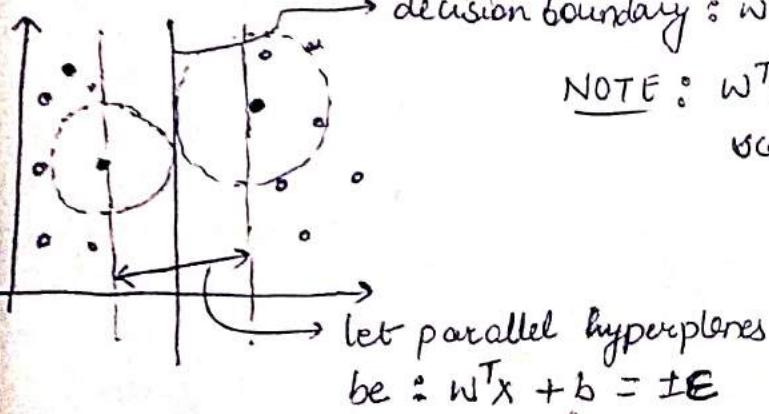


- Maximizing margin improves generalization  
 →  $h$  can be made independent of the dimensionality :  $d$ .

### Formalizing the margin

decision boundary :  $w^T x + b = 0$ .

NOTE :  $w^T x_i + b$  is dependent on the scale of  $X$  and  $w$ .



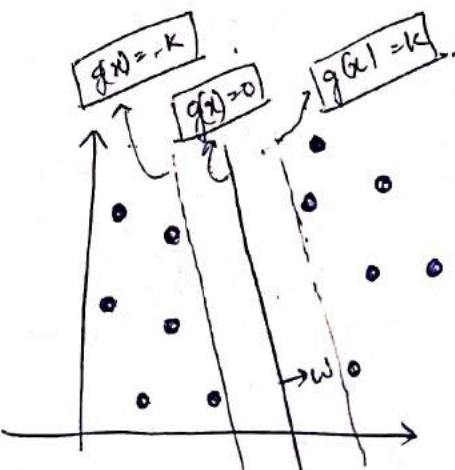
### Formulation

→ let  $g(x) = w^T x + b$ .

→ We want to maximize such that

$$-w^T x_i + b \geq k \text{ for } d_i = 1.$$

$$-w^T x_i + b \leq -k \text{ for } d_i = -1.$$



→ Values of  $g(x)$  depends on  $\|w\|$ :

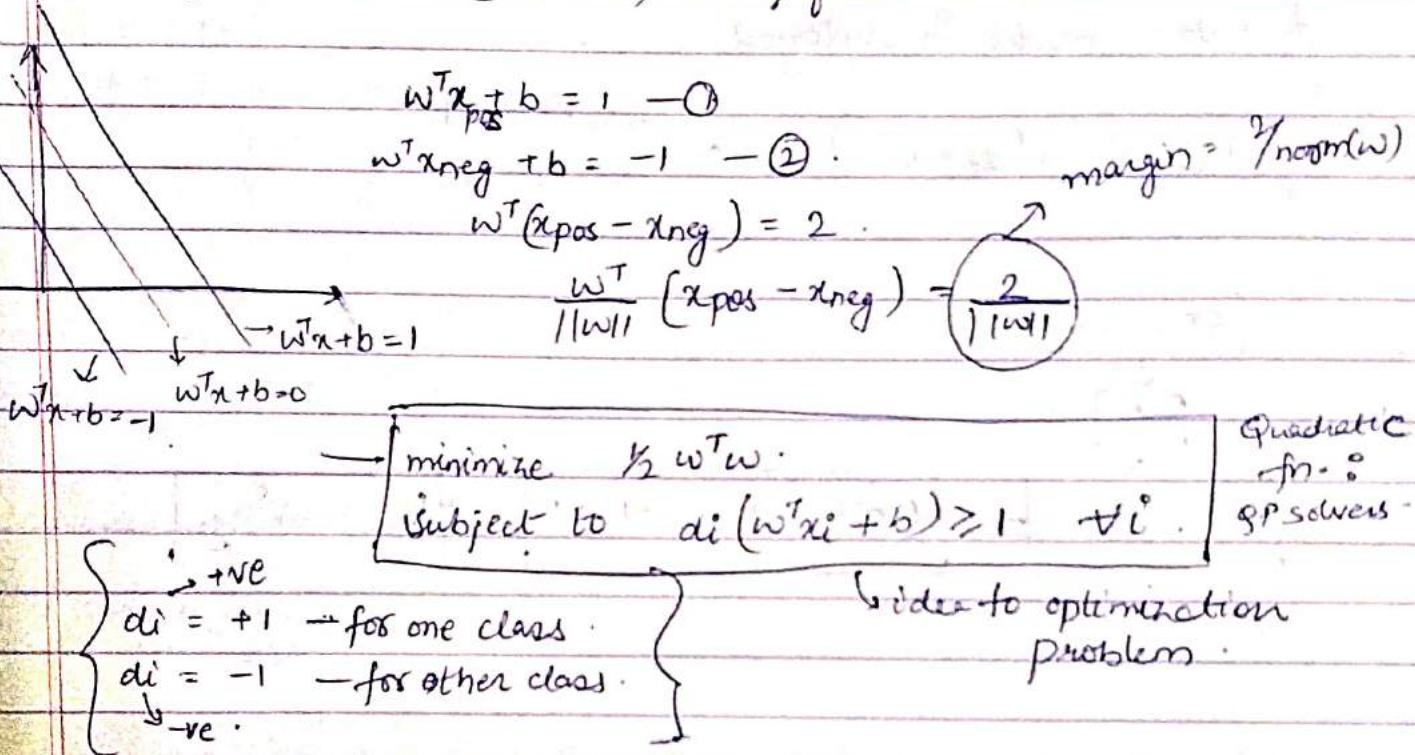
1. Keep  $\|w\| = 1$ , and maximize  $g(x)$ , or.

2. Let  $g(x) \geq 1$  and maximize  $\|w\|$ .

→ We use two approaches and formulate the problem as:

1. Minimize:  $\frac{1}{2} w^T w$ .

2. Subject to:  $d_i(w^T x_i + b) \geq 1$ , for  $i = 1, \dots, N$ .



SMAI (04/02/2020)

$$\rightarrow J(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i d_i (w^T x_i + b) + \sum d_i^2$$

S.t.  $d_i \geq 0 \quad \forall i$

Lagrangian form

$$w^* = \sum_{i=1}^N \alpha_i d_i x_i$$

$$\frac{\partial J}{\partial w} = 0$$

$$\sum_{i=1}^N \alpha_i d_i = 0$$

$$\frac{\partial J}{\partial b} = 0$$

Minimize w.r.t.  
 $a, b$  and max.  
w.r.t.  $d$ .

$$\alpha_i (\text{di}(\mathbf{w}^T \mathbf{x}_i + b^*) - 1) = 0$$

$$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j \text{di}(\mathbf{w}^T \mathbf{x}_i + b^*) \text{di}(\mathbf{w}^T \mathbf{x}_j + b^*)$$

$\alpha_i$

$$\mathbf{w}^* = \sum \alpha_i \mathbf{d}_i \mathbf{x}_i ; b^* = 1 - \mathbf{w}^T \mathbf{x}_{S^*}$$

→ Now model is deployed.

$$\left\{ \begin{array}{l} \text{sign } (\mathbf{w}^T \mathbf{x}_i + b^*) \rightarrow \\ \quad \text{+ve : class 1} \\ \quad \text{-ve : class 2.} \end{array} \right\}$$

QP Solver

$\downarrow$   
 $\alpha_i$   
 $\downarrow$

$$\mathbf{w}_0 = \sum_{i=1}^N \alpha_i \mathbf{d}_i \mathbf{x}_i \rightarrow \text{di}[\text{di}(\mathbf{w}_0^T \mathbf{x}_i + b_0) - 1] = 0 \rightarrow b_0 = 1 - \mathbf{w}_0^T \mathbf{x}_{S^*}$$

where  $\alpha$  is non zero in SV.

How to modify current program statement to support noise?

$$\begin{aligned} & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t. } & \text{di}(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 + \varepsilon_i \end{aligned}$$

$\varepsilon_i$  → Slack Variable.

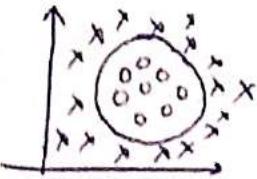
modified constraint :

$$\text{di}(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \varepsilon_i \quad \forall i$$

## Non-Separable Data

1. Noisy data / Bad features

2. Non-linear boundary.



## SVM with noisy data

$$\text{Minimize : } \phi(w) = \frac{1}{2} w^T w.$$

$$\text{Subject to : } d_i(w^T x_i + b) \geq 1 - \xi_i$$

↓  
Introduce slack variables  $\xi_i \geq 0$ ,

$$\text{Minimize : } \phi(w) = \frac{1}{2} w^T w.$$

$$\text{Subject to : } d_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\text{Minimize : } \phi(w) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i$$

$$\text{Subject to : } d_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad \forall i$$

## Dual form with Slack Variables

$$Q(x) = \sum \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N x_i \alpha_j d_j d_i x_j^T x_i$$

$$\text{Subject to } \alpha_i \geq 0 \text{ and } \sum_{i=1}^N \alpha_i d_i = 0.$$

→ Note that neither the slack variables, nor their lagrange multipliers appear in dual.

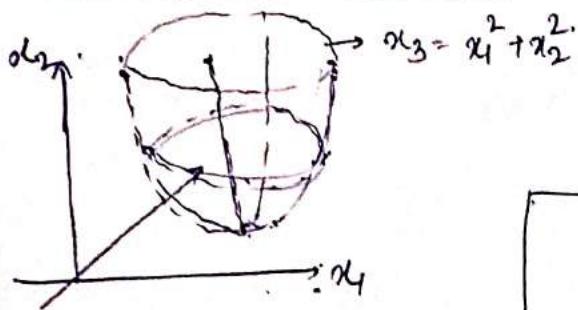
→ The only change is the additional constraint on  $\alpha_i$ .

→ The parameter  $C$  controls the relative weight between training error and the VC dimension.

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i \alpha_j d_j \alpha_i x_i^T x_j$$

$$\text{Subject to } 0 \leq \alpha_i \leq C \quad \forall i \text{ and } \sum_{i=1}^N \alpha_i d_i = 0.$$

## Non-linear Boundaries



$$\phi(x) = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{pmatrix}$$

Training -

### SVM post mapping

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \phi(x_i) \cdot \phi(x_j)$$

Subject to  $0 \leq \alpha_i \leq C$   $\forall i$  and  $\sum_{i=1}^N \alpha_i d_i = 0$ .

$$\text{Label} = \text{sign} (\omega_0 \cdot \phi(x_{\text{test}}) + b_0)$$

$$\omega_0 = \sum_{i=1}^N \alpha_i d_i \phi(x_i)$$

$$\therefore \text{label} = \text{sign} \left( \sum_{i=1}^N \alpha_i d_i \phi(x_i) \cdot \phi(x_{\text{test}}) + b_0 \right)$$

Testing -  
for testing  
also we have  
to transform point to 3-D and then check.

### SVM post Mapping

- Data vectors occur only as dot products in SVM - learning & testing.
- If we can find a fn.  $K(x, Y)$ , which is equivalent to  $\phi(x) \cdot \phi(Y)$ , we can avoid

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(x_i, x_j)$$

$$\text{label} = \text{sign} \left( \sum_{i=1}^N (\alpha_i d_i K(x_i, x_{\text{test}})) + b_0 \right).$$

### What do we gain by using $K$ ?

original space : 2-dimensional

$$\text{let } K(x, Y) = (x \cdot Y)^3 = (x_1 y_1 + x_2 y_2)^3$$

$$\text{let } \phi(x) = \phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} x_1 x_1 \\ x_1 x_2 \\ x_2 x_1 \\ x_2 x_2 \end{pmatrix}$$

$$\text{let } K(x, y) = \phi(x) \cdot \phi(y) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2 x_1 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} y_1^2 \\ y_1 y_2 \\ y_2 y_1 \\ y_2^2 \end{bmatrix}$$

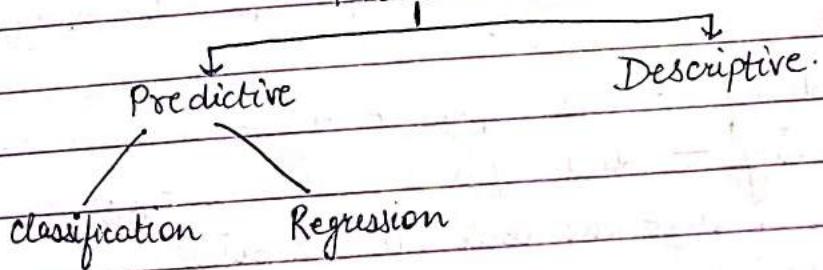
$$= x_1^2 y_1^2 + 2x_1 x_2 y_1 y_2 + x_2^2 y_2^2$$

$$= (x_1 y_1 + x_2 y_2)^2$$

$$= (x \cdot y)^2$$

## Unsupervised learning - Clustering (15/02/2020)

→ clustering news of same type (topic) - Google news.  
 ML Tasks



ML Tasks :: Descriptive -

- Study / Exploit the 'structure' of data.
- clustering
- Dimensionality Reduction
- Density Estimation
- Also called as 'unsupervised learning'
- Input data without paired output

→ cluster images - sea, sunset, flowers etc.

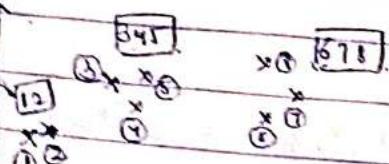
- Determine groups of people in image
- based on clothing styles } features.
- gender, age etc.
- Determine moving objects in videos.

What is clustering?

- high intra-group similarity (within members of a cluster)
- low inter-group similarity (across clusters)

### Agglomerative Clustering

considered  
as  
single  
point



	1	2	3	4	5	6	7	8
1	0							
2		0						
3			0					
4				0				
5					0			
6						0		
7							0	
8								0

Symmetric matrix.

→ we calculate distance of every point with every other point

→ we find closest point; make 1 & 2 as single point; 1 & 2 can be defined as mean of 1 & 2.

→ Average link clustering — we represent cluster as mean of points.

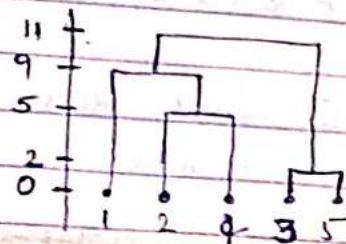
→ Single link clustering — dist. b/w 2 clusters is equal to dist b/w their closest point.

→ Complete link clustering — dist. b/w 2 clusters is the dist b/w their farthest point.

3 ways of agglomerative clustering based on distance computation.

→ Start with clusters equal to no. of points and then keep merging them.

	1	2	3	4	5
1	0				
2		0			
3			0		
4				0	
5					0



Dendrogram.

	3	5	1	2	4
3	0				
5		0			
1			0		
2				0	
4					0

$$J = \sum_{k=1}^K \sum_{i=1}^{n_k} \|x_{ki} - \mu_k\|^2 \quad \leftarrow \text{to minimize in case clustering.}$$

### Beehive clustering

- decide threshold = 0.3 (suppose).
- when new pt. comes, compute distance from  $c_i$  & either give it to same cluster or form a new cluster.
- c. Done for every new point. Compute distance from each cluster & work around.
- faster computationally.
- Sir ki Khoj.
- ⇒ In clustering, distance metric to be used is also crucial. (cosine, euclidean, manhattan).
- clustering and then KNN.  $\Rightarrow$  very effective.

### K-Means Clustering

- Approach 1: Assume that we (somehow) know no. of clusters.
- # of clusters =  $k$  (often a 'guess').
  - How to represent a cluster?
    - Cluster members
    - A suitable statistic of the cluster numbers.

K-Means Clustering : ① Initialization  $\rightarrow$  decide  $k$ , and initialize  $k$  centers (randomly).

→ Take a point, see which chosen point is it closest to and assign it to that cluster.

- ② Iteration 1 — assign all objects to the nearest center; move a center to mean of its members
- ③ Iteration 2 — After moving centers, reassign the objects to nearest centers. Move a center to the mean of its new members.

Reassign and move centers; until no objects changed membership.

Time complexity =  $O(k.n)$ .  $\leftarrow$  practically faster; k means ~~can't~~ always converge but agg. converges theoretically.

{for agglomerative clustering — TC =  $O(n^2)$ }.

$$\{x^{(1)}, \dots, x^{(m)}\} ; x^{(i)} \in \mathbb{R}^n$$

The k-means clustering algo is as follows:

① Initialize cluster centroids  $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$  randomly.

② Repeat until convergence.

for every set  $i$ , set

$$c^{(i)} = \arg \min \|x^{(i)} - \mu_j\|^2 \rightarrow \text{assignment step: assign each data pt to closest cluster.}$$

for each  $j$ , set

$$\mu_j = \frac{\sum_{i=1}^m \mathbf{1}\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)} = j\}} \quad \begin{matrix} \nearrow \text{Refitting step:} \\ \text{Move each cluster center} \\ \text{to the center of the data} \\ \text{assigned to it.} \end{matrix}$$

?

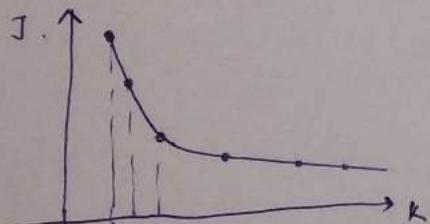
### Why k-means Works?

→ high intra cluster similarity.

→ K-means optimizes  $J = \sum_{k=1}^K \left[ \sum_{i=1}^{n_k} \|x_{ki} - \mu_k\|^2 \right]$

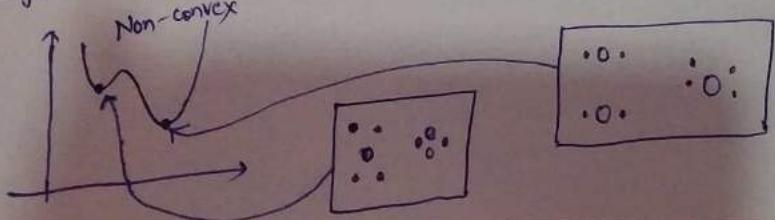
closed clustering - this term will be decreased (dist from mean).

### How to decide K?



Plot for diff values of  $K$  & choose  $K$  where elbow appears generally.

⇒ Objective function for k-means is non-convex



⇒ for  $K$ , try with diff initialisations & pick  $K$  which gives lowest  $J$ .

SMA1 (01/03/2020)

### Initialization

- ① Sample a small set of points  
 ↓  
 Agglomerative clustering

### K-means ++

- ① Randomly choose one of the observation to be a cluster center
- ② For each observation  $x$ , determine  $d(x)$ , where  $d(x)$  denotes the minimal distance from  $x$  to a cluster center.
- ③ choose next cluster center;  $(d(x))^2$
- ④ Repeat 2 & 3 till you have  $k$  clusters

### SVM Kernels.

→ experiment common kernels & find what works best

ML - A probabilistic perspective. (Generative Modelling)

- Mean, variance, co-variance, expectation etc.
- Gaussian dist
- Sampling from a dist

$$\text{Cov.} = \sum (x - \mu_x)(y - \mu_y)$$

	$x$	$y$
$x$	$\sigma_x^2$	$\sigma_{xy}^2$
$y$	$\sigma_{xy}$	$\sigma_y^2$

$$\text{Expectation} = \sum_i x_i P(x=x_i) \\ E(x)$$

$$E[(x-\mu)^2] = \sum (x_i - \mu)^2 \cdot P(x_i)$$

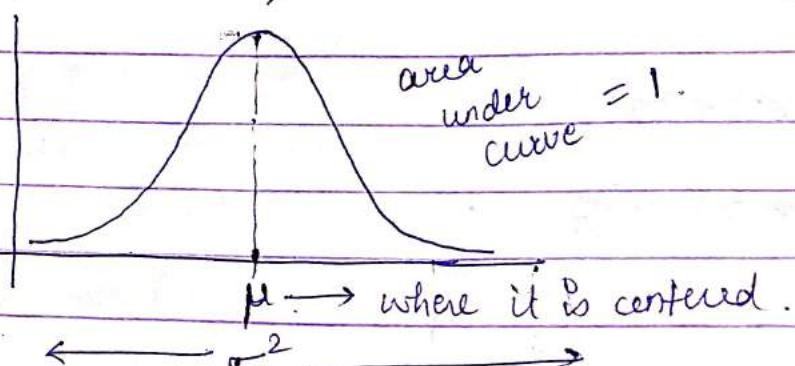
$$V(x) = E[x^2] - (E[x])^2$$

### Gaussian distribution (Normal)

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

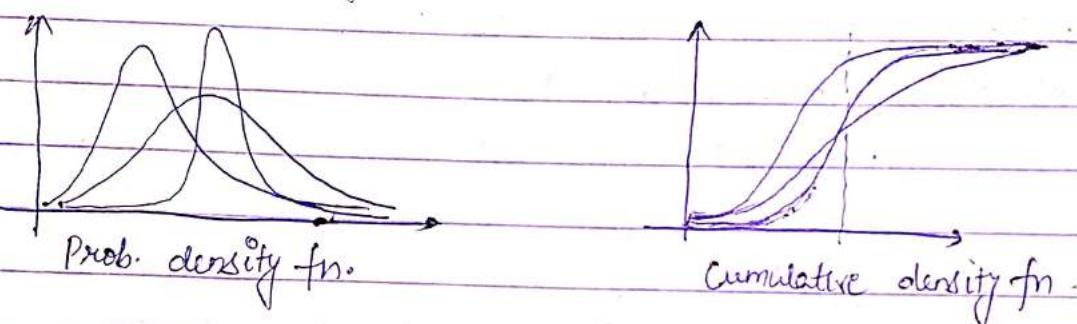
} normal  
Gaussian

shows how likely a data can appear.



↳ how much spread.

II Sample from a gaussian distribution →



- Generative vs discriminative classification
- density estimation - anomaly detection
- Gaussian distribution (univariate & multivariate)
- Generative classification (Naive Bayes).

SVM - discriminative classification, we discriminate on the basis of hyperplane. (divide).

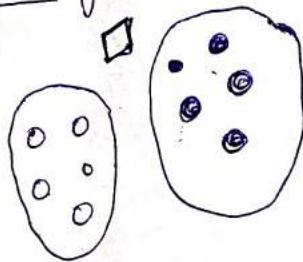
Generative - fitting a distribution to each class.

Given a ~~distribution~~ <sup>sample</sup> class, we try to find which class discrimination it belongs

Not discriminating, doing generative modelling.

- fit some density fn. to each class.

## Density Estimation



Compute  $\begin{cases} p(x|y_1) \\ p(x|y_2) \end{cases}$

→ Given a point, how do I fit a fn. to it.

2-D Gaussian.

- ~~Hence~~ Most of the processes we see in real world, is actually a cumulation of many random events.
- As you take many processes and sum their effect, you tend towards a gaussian fn.

## Anomaly Detection

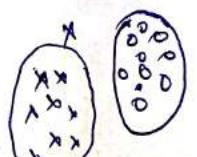
Is everything fine with my aircraft? → 2 class classification prob.  
If less negative data present, calculate prob. based on some values and then set threshold like if pdf values is less than 0.1, then perform manual check.

→ Multivariate Gaussian Distribution.

$$f(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu))}.$$

## SMAI (04/03)

likelihood.



$p(x|y_1) \quad p(x|y_2)$  — we can calculate this.

~~discriminate~~ discriminative classification — we <sup>are learning</sup> have to discriminate b/w the two

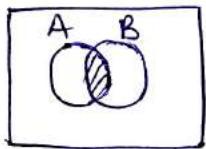
Generative classification — model each of the class & based on model, likelihood. Compute

$$p(y_i|x) \quad p(y_i|x) \longrightarrow \text{from pdf -}$$

decision rule -  $\operatorname{argmax}_i p(y_i|x)$

How to compute  $p(x|y)$  - use Bayes Rule.

$$p(y|x) = \frac{p(x|y) \cdot p(y)}{p(x)}$$



$$P(A|B) = \frac{P(A \cap B)}{P(B)} ; \quad P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

$$\Rightarrow \text{Bayes Rule} - P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

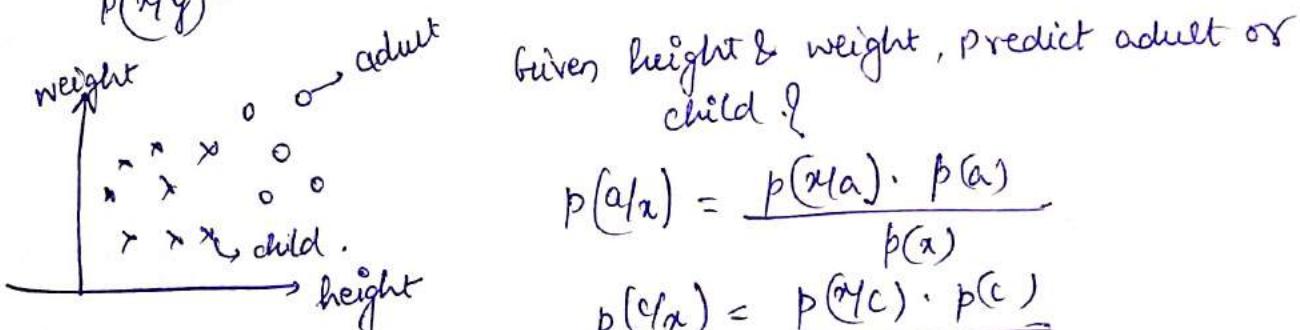
$\Rightarrow$  decision - comb. of what we observe & what priors we have.

### Naive Bayes

$$P(x_1, x_2, \dots, x_d | y) = \prod_{i=1}^d p(x_i|y) \quad \xrightarrow{\text{computing likelihood of each feature}}$$

$\rightarrow$  it fails when features depend on each other.

$$p(x|y)$$



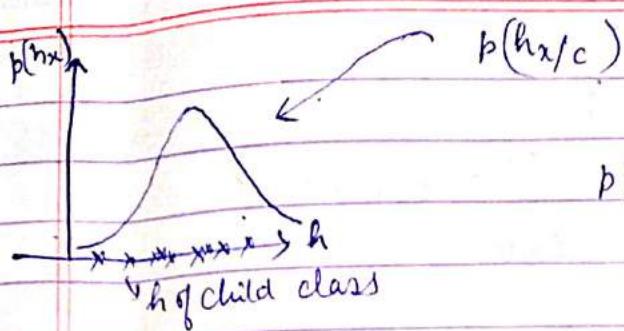
$$p(a|x) = \frac{p(x|a) \cdot p(a)}{p(x)}$$

$$p(c|x) = \frac{p(x|c) \cdot p(c)}{p(x)}$$

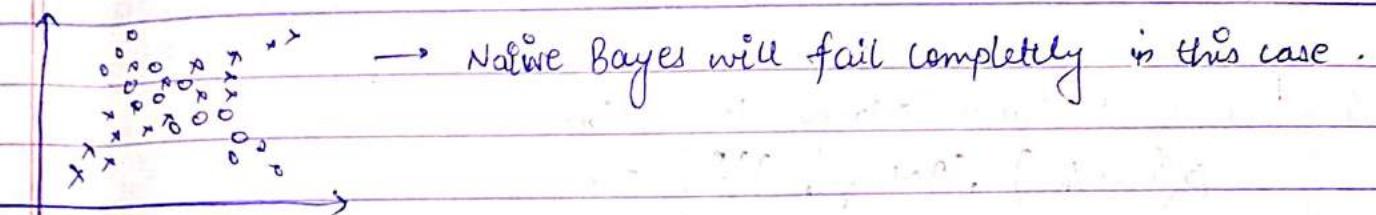
$$p(a) = 0.25 ; p(c) = 0.75$$

$$p(x|c) = p(h_x|c) \cdot p(w_x|c) ; p(x|a) = p(h_x|a) \cdot p(w_x|a)$$

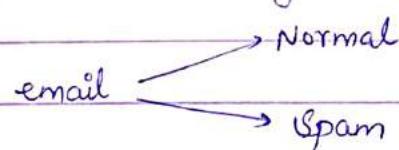
$$x = [h_x, w_x]$$



$$p(c/x) = \frac{p(x/c) \cdot p(c)}{p(x/a) \cdot p(a) + p(x/c) \cdot p(c)}$$



⇒ Multinomial Naïve Bayes — Text Classification:



[ Once you model a class  
you can sample that  
class.  
No CV, NO req. ]

[ dear, friend, lunch, money ]

$$\text{Normal} - P(\text{dear}/N) = 0.47.$$

$$P(\text{friend}/N) = 0.29$$

$$P(N) = 0.67.$$

$$P(\text{lunch}/N) = 0.18$$

$$P(\text{money}/N) = 0.06$$

$$\text{Spam} - P(\text{dear}/S) = 0.29$$

$$P(\text{friend}/S) = 0.14$$

$$P(S) = 0.33$$

$$P(\text{lunch}/S) = 0$$

$$P(\text{money}/S) = 0.57$$

Based  
on  
training  
data

E-mail  $\Rightarrow$  "Dear Friend"  $p(N/x) \cdot P(S/x)$

$$= p(x/N) \cdot P(N)$$

$$= P(\text{dear}/N) \cdot P(\text{friend}/N) \cdot P(N)$$

$$= 0.47 \times 0.29 \times 0.67 = 0.09 \propto P(N/\text{dear friend})$$

$$\begin{aligned}
 &= p(x/s) \cdot p(s) \\
 &= p(\text{dear}/s) \cdot p(\text{friend}/s) \cdot p(s) \\
 &= 0.29 \times 0.14 \times 0.33 = 0.01 \propto P(S/\text{Dear friend})
 \end{aligned}$$

→ So, we classify it to Normal class.

→ Use words from vocabulary only; don't use is, am, are for classification.

Q. E-mail : Lunch Money Money Money.

$$p(\text{lunch}/n) \cdot (P(\text{Money}/n))^3 \cdot P(n)$$

$$= 0.00002 \cdot (\text{almost}) \propto P(n/x).$$

$$p(\text{lunch}/s) \cdot (P(\text{Money}/s))^3 \cdot P(s)$$

$$= 0. \propto P(s/x)$$

→ Should be classified to Spam but due to lunch being 0, didn't happen.

→ So, we assign min. value (Suppose 0.01 is added everywhere).

$$\text{Now, } (0.01) \cdot (0.58)^3 > 0.0002 \rightarrow \text{Hence Spam.}$$

Gaussian.

# Can you apply Naive Bayes to Textual data?

No, bcz. there isn't continuous variable, we have Categorical data.

Not used together.

Naive Bayes — Works very good on textual data.

decision trees — textual → — data.

challenge - missing data - add avg.  
 → prob. may not add upto 1.

### Text classification

- Naive Bayes ✓
- Bag of Words

→ SVM not used - bcz numerical data; fixed dim. representation.

[Stop word removal - high freq words - is, am, an, the]

$$\begin{array}{l}
 \text{love} \\
 \text{cinema} \\
 \text{genius} \\
 \text{happy} \\
 \text{↑ movie} \\
 \text{vocab.}
 \end{array}
 \left[ \begin{array}{c} 2 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 2 \end{array} \right] / \text{words.} = \left[ \begin{array}{c} 2/1 \\ 1/1 \\ 1/1 \\ 0 \\ \vdots \\ 2/1 \end{array} \right] \leftarrow \begin{array}{l} \text{representation} \\ \text{of} \\ \text{textual} \\ \text{data.} \end{array}$$

$$\# \text{tf-idf} = \text{tf}(t_1 D) \cdot \text{idf}(t_1 D).$$

$$\text{tf } t_1 D = \log(1 + \text{freq}(t_1 D))$$

$$\text{idf } (t_1 D) = \log \left( \frac{N}{\text{no. of doc. having term } D} \right).$$

$$\begin{array}{l}
 w_0 \\
 w_1 \\
 \vdots
 \end{array}
 \left[ \begin{array}{c} \text{tf}(w_0) \cdot \text{idf}(w_0) \\ \text{tf}(w_1) \cdot \text{idf}(w_1) \\ \vdots \end{array} \right]$$

→ words which occur rarely, gets more tf-idf value in this method.

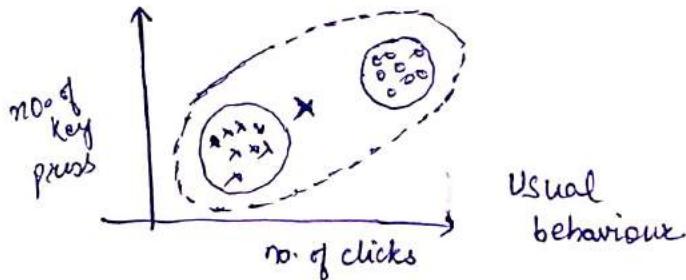
→ bag of words is just normalized freq.

→ both don't consider order of words.

$$\begin{array}{l}
 \text{bigram} \\
 \text{bigram} \\
 \vdots
 \end{array}$$

## Gaussian Mixture Models

→ Anomaly Detection



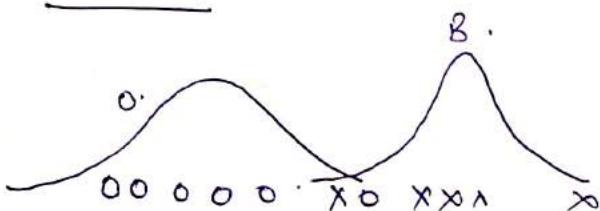
K-means

- hard assignment of clusters.
- what if clusters are overlapping?

GMM

- clusters are modeled as gaussians (not just by mean).
  - ↳ trained using EM algo - generative model..
  - ↳ you can sample from it.

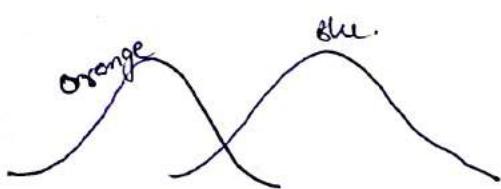
Scenario 1



→ if I know the source, I can estimate the (gaussians) parameters.

Scenario 2 - parameters of gaussian are given

→ I know the gaussians



$$P(b|x_i) = \frac{P(x_i|b) \cdot P(b)}{P(x_i|b) \cdot P(b) + P(x_i|a) \cdot P(a)}$$

$P(x_i|b)$  = substitute in formula.

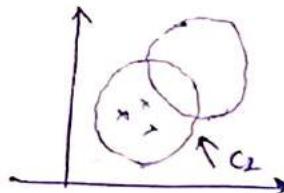
these prob. can be computed in this scenario.

→ real world data might not be a single gaussian.

That's why need GMM.

$$f(x) = \sum_c \pi_c N(x; \mu_c, \sigma_c)$$

$\{\sum \pi_c = 1\}$ . weighted sum of gaussian.



- GMM assumes how many clusters. Suppose  $K=2$ .
  - EM Algo - [Expectation Maximization Algo]
  - E → ① Start with two randomly placed Gaussians  $(\mu_a, \sigma_a^2), (\mu_b, \sigma_b^2)$
  - E → ② For each point  $x_i$ , compute  $P(b_i/x_i) \rightarrow$  does it look like it came from b.  
 $P(a_i/x_i) \rightarrow$  " " " " from a.
  - M → ③ Adjust  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$  to fit points assigned to them.
  - Here soft clustering, not hard like k-means.

$$\text{Expectation} \rightarrow P(x_i|b) = \frac{1}{\sqrt{2\pi\sigma_b^2}} e^{-\frac{(x-\mu_b)^2}{2\sigma_b^2}}$$

Step

$$b_i = P(b|x_i) \rightarrow \frac{P(x_i|b) P(b)}{P(x_i|b) P(b) + P(x_i|a) P(a)}$$

$$x_i = 1 - b_i$$

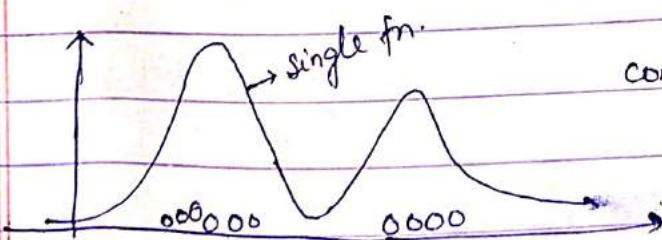
$$\text{Maximization} \rightarrow$$

Step  $\mu_b = \frac{b_1 x_1 + b_2 x_2 - b_n x_n}{b_1 + b_2 - b_n} ; \sigma_b^2 = \frac{b_1 (x_1 - \mu_b)^2 + b_2 (x_2 - \mu_b)^2 + \dots + b_n (x_n - \mu_b)^2}{b_1 + b_2 + \dots + b_n}$

$$\mu_a = \checkmark \quad ; \quad r_a^2 = \checkmark$$

$$P(a) = \pi_a = \frac{a_1 + a_2 - \dots - a_n}{N}$$

$$\pi_a = 1 - \pi_b$$



convex combination of a set of gaussian.

$\rightarrow$  C gaussian -  $\underbrace{\pi_1, \pi_2, \dots, \pi_K}_{0.1 \quad 0.01 \quad 0.4 \quad \dots \quad 0.1}$

pick (find) a gaussian, know its parameters. So,

Sample point from that gaussian.

## 2-way sampling

## SMAI (15/03).

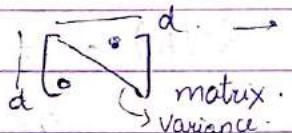
Principal Component Analysis (PCA) - for dimensionality reduction  
 feature dimension - 1000 dimensional vectors as input.  
 ✓  $\rightarrow$  98%.

$\rightarrow$  unwanted data goes away }  
 $\rightarrow$  more structured data. }  
 ✓  $\rightarrow$  97%.

$\rightarrow$  98% accuracy.

bcz lot of unwanted columns are reg.

Algo 1 - covariance matrix



d.  $\rightarrow$  look at high values & keep one variable among highly correlated.

Algo 2 - classification - 1d, either add feature and then if acc. increase with best acc. then keep them.

Algo 3 - dropping from 100 d.

PCA - defines each dimension as linear weighted sum of other dimension  
 ↳ compression

↳ PCA prob - How to retrieve this?  
 prob - to find axis on which to project.

$\rightarrow$  if data is projected on  $u_1$  - data reduced to 1-D.  
 ↳ variance is maximized.

#1 Maximizing Variance - leads to maximizing info.

$$\{ \text{entropy} = -\sum p_i \log p_i \}$$

↳ measure of some sort of randomness ; variation

$\rightarrow$  small direction, less variance, ~~better~~ better prediction.

$$\boxed{\text{Covariance} \rightarrow \frac{1}{N-1} \sum_{i=0}^N (x_i - \mu_x)(y_i - \mu_y)}$$

↳ each pair one by one.

$$x = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix}_{d \times n}. \quad x_i = \begin{bmatrix} ? \\ \vdots \\ ? \end{bmatrix}_d$$

vector

$$\text{cov}(x_1 x_2) = \frac{1}{N-1} \sum_{i=1}^N (x_1^i - \mu_{x_1})(x_2^i - \mu_{x_2})$$

$$\text{cov}_i = \frac{1}{N-1} \sum_{i=0}^N (x_i^i - \mu_0)(x_i^i - \mu_1)$$

$$\left\{ \frac{1}{N} (x_1 + x_2 + \cdots + x_N) \right\}$$

μ vector.

$$(x - \mu) \rightarrow [x_1 - \mu \quad x_2 - \mu \quad \cdots \quad x_n - \mu]_{d \times n}.$$

$$\frac{(x - \mu)}{d \times n} \frac{(x - \mu)^T}{n \times d} \rightarrow d \times d \rightarrow \text{covariance matrix.}$$

$$\begin{bmatrix} v_{00} & v_{01} & v_{02} & \cdots \\ v_{10} & v_{11} & \cdots & \cdots \\ \vdots & \vdots & \ddots & \ddots \end{bmatrix}$$

$$Y = \frac{\mu^T}{d \times n} X_{d \times n} \rightarrow 1 \times n.$$

$$Y = [y_1 \ y_2 \ \cdots \ y_n]$$

$$y_i \in \mathbb{R}.$$

$$\boxed{\max_u u^T x - \text{var}(u^T x)} \quad \text{(optimization prob.)} \rightarrow E(u^T x - E[u^T x]) - (u^T x - E[u^T x]^T).$$

$\left\{ \begin{array}{l} E(x) = \mu \\ E[u^T x] = u^T \mu \end{array} \right.$

scalar times mean

$$\rightarrow u_1^T E[(x - \mu)(x - \mu)^T] u_1$$

$$\boxed{\text{Var}(u_1^T x) = u_1^T S u_1.} \quad \leftarrow \text{maximize this.}$$

$$\Rightarrow \text{Maximize } u_1^T S u_1 \text{ wrt } u_1 \rightarrow \text{quadratic fn.} \quad \boxed{U} \rightarrow \text{no maxima.}$$

S.t.  $u_1^T u_1 = 1.$

$$\boxed{\begin{array}{l} \max_{u_1} u_1^T S u_1 \\ \text{s.t. } u_1^T u_1 = 1 \end{array}} \rightarrow \text{optimization prob. for PCA.}$$

$$J(U) = \frac{U_1^T S U_1}{\text{scalar}} - \lambda (U_1^T U_1 - 1)$$

$\stackrel{\text{scalar}}{\cancel{\frac{1 \times d}{d \times d} \frac{d \times d}{d \times 1}}} = 1 \times 1$

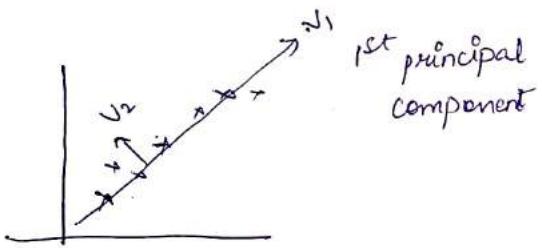
diff. scalar w.r.t. vector  
 we get vector if  
 same dimension.

$$J'(U) = 0 \Rightarrow 2S U_1 - 2\lambda U_1 = 0.$$

$$S U_1 = \lambda U_1 \rightarrow \text{KEY RESULT:}$$

eigen vector.

$$\left\{ \begin{array}{l} \max. J(U) = U_1^T \lambda U_1 \\ \max. U_1 = \lambda \end{array} \right\}.$$



$\lambda_1, v_1$  → direction vector.  
 $v_2 \rightarrow$  orthogonal.

convert 2D to 1D.

pick larger eigen vector i.e. with  
 larger  $\lambda$ . ( $\lambda_1 > \lambda_2$  suppose)

$$\left[ \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right] \xrightarrow{\text{2d.}} \underline{\underline{\text{cov}}}.$$

$$\left\{ \begin{array}{l} (\lambda_1, v_1) \\ \lambda_2 v_2 \end{array} \right\} \text{ 2 eigen vectors.}$$

$$p = [p_1, p_2, \dots, p_d]$$

$$p_1^2 + p_2^2 + \dots + p_d^2 = 1 \rightarrow \text{unit vector.}$$

$$p^T p = [p_1, p_2, \dots, p_d] \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_d \end{bmatrix}$$

$$p^T p = 1 \rightarrow \text{direction.}$$

{ PCA — linear method }  
 kernel PCA exist.

## SVD (Similar Value decomposition)

$$X = U \Sigma V^T$$

eigen vector  
 of  $XX^T$       eigen vectors  
 $X^T X$ .  
 diagonal matrix  
 $\begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$

$\left. \begin{array}{l} \tilde{X} = (X - M) \\ d \times d \times N \\ X^T \tilde{X} \end{array} \right\} \begin{array}{l} \text{covariance} \\ \text{matrix} \end{array}$

→ when we do SVD on  $X$

$$X = [x_1 \dots x_d]$$

$$\tilde{X} = [x_1 - \mu \quad x_2 - \mu \quad \dots \quad x_n - \mu]_{d \times n} \rightarrow \text{SVD}$$

Procedure for PCA.

$d$ -dim. data

→ SVD →  $d$  eigen values.

$\left\{ \begin{array}{l} \lambda_1 \lambda_2 \dots \lambda_d \\ u_1 u_2 \dots u_d \end{array} \right\} \rightarrow d$  eigen vector.

If 1-D — multiply by  $\lambda_1, u_1$ .

2-D — "  $\lambda_1 u_1 + \lambda_2 u_2$ .

for  $p$ -D data — multiply by 1<sup>st</sup>  $p$  vectors.

How many dimension should we pick?

$$r_p = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_p}{\sum \lambda_i} = 0.99$$

take  $p$  in such way.

$$d = 1500$$

$$p = 40.$$

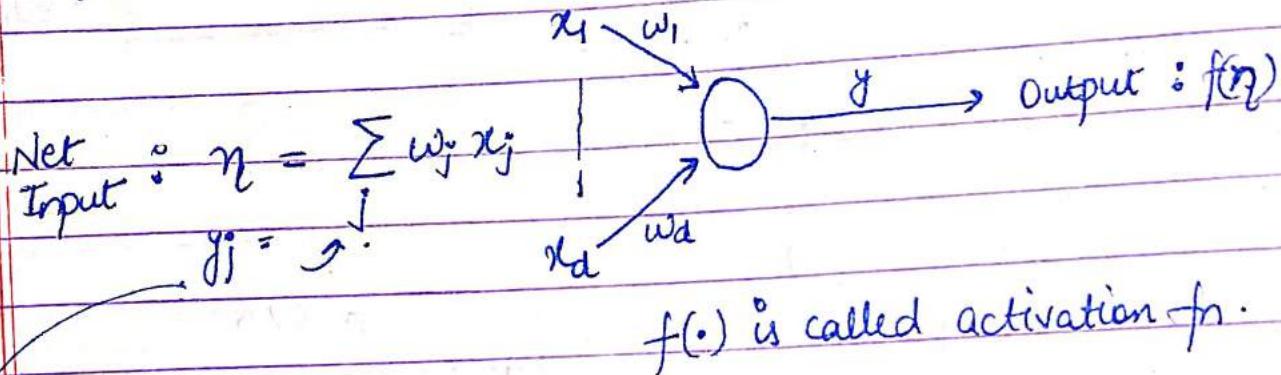
# SMAI (18/03)

Neural Networks.

→ Artificial Neural N/W.

→ parallel distributed info processor. —  
{Slides}

## Single Neuron Model

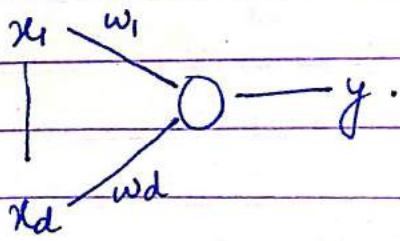


$$\{ L = \sum_j (y_j - \hat{y}_j)^2 \} \rightarrow \text{loss function.}$$

$$w_{t+1} \rightarrow w_t - \lambda \frac{\partial L}{\partial w}.$$

## Perception

$$y = \sum_i w_i x_i = w^T x.$$



$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

2 class classification

$$\left\{ \begin{array}{l} w^T x \geq t \rightarrow \text{class 1} \\ w^T x < t \rightarrow \text{class 2} \end{array} \right\}$$

$$y = w^T x + b \quad (\rightarrow \text{bias})$$

$\text{sign}(w^T x + b)$   $\rightarrow$  classification

Goal - to learn how to calculate weights.

$$x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$w^T x = w_0 + \sum_{i=1}^d w_i x_i$

$\rightarrow$  Rosenblatt (1962)

$$\Delta w_k = w_{k+1} - w_k$$

$$\Delta w_k = 0 \quad \left\{ \begin{array}{l} \text{if } w_k^T x > 0 \text{ & } y(k) = 1, \\ \text{or } w_k^T x < 0 \text{ & } y(k) = 0 \end{array} \right.$$

$$\Delta w_k = \begin{cases} x(k) & \text{if } w_k^T x \leq 0 \text{ & } y(k) = 1 \\ -x(k) & \text{if } w_k^T x > 0 \text{ & } y(k) = 0. \end{cases}$$

Prediction  $\rightarrow \text{sign}(w^T x)$

made a mistake  $\rightarrow y(k) = 1, \quad \boxed{w_k^T x \leq 0}$ .

$$w_{k+1} = w_k + x_k x_i$$

$$\boxed{w_k^T x}$$

old prediction,  
made an error

$$\longrightarrow \boxed{w_{k+1}^T x}$$

new prediction

$$w_{k+1}^T x = (w_k + x_k)^T x$$

$$\boxed{w_{k+1}^T x = w_k^T x + (x_k)^2}.$$

$\hookrightarrow$  prev. prediction + a small positive value.

## Perception learning Algo — slide.

→ In finite no. of iterations, perceptron algo finds correct line.

$$w^T x_1 + b = 0$$

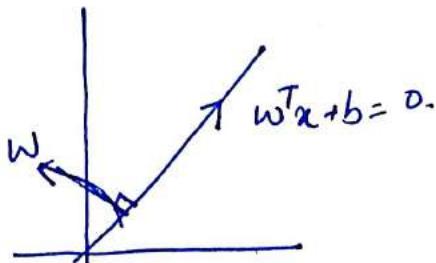
$$w^T x_2 + b = 0$$

$$b + w^T x_1 \neq w^T x_2 + b$$

$$w^T(x_1 - x_2) = 0$$

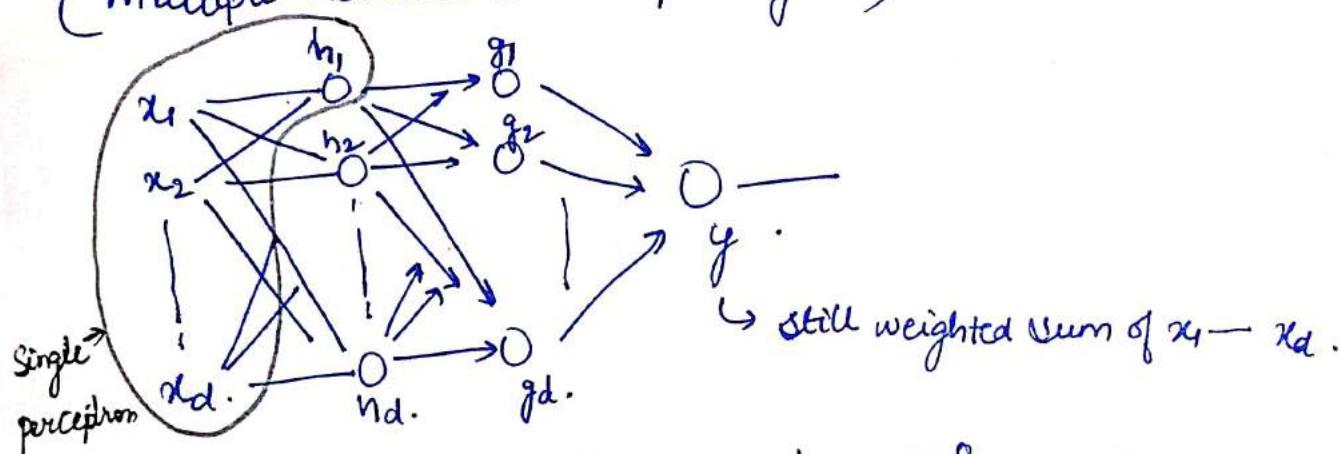
↓

Vector • Vector:



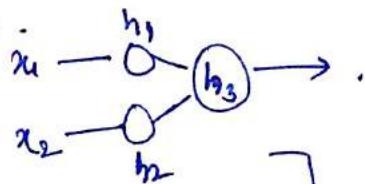
## Multi-layer Perceptron [— learning a linear fn. only]

(Multiple neurons in multiple layers)



still weighted sum of  $x_1 - x_d$ .

⇒ Adding multiple layers of neurons is a chain.



$$h_1 = w_{11}x_1 + w_{12}x_2$$

$$h_2 = w_{21}x_1 + w_{22}x_2$$

$$h_3 = w'_{11}h_1 + w'_{12}h_2$$

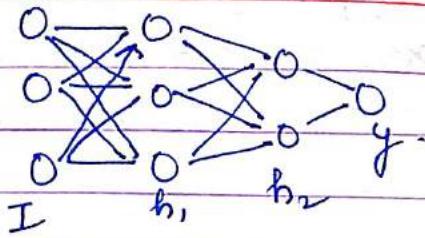
$$h_3 = w'_{11}(w_{11}x_1 + w_{12}x_2) + w'_{12}(w_{21}x_1 + w_{22}x_2)$$

$$h_3 = (w'_{11}w_{11} + w'_{12}w_{21})x_1 + (w'_{11}w_{12} + w'_{12}w_{22})x_2$$

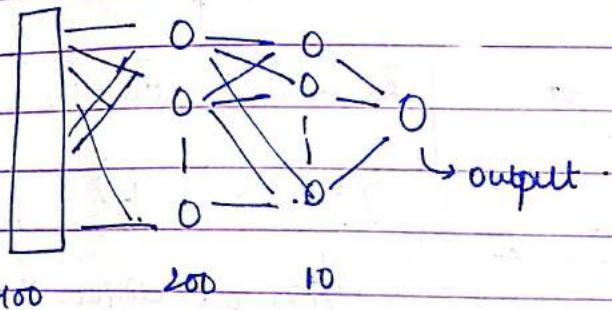
↪  $h_3$  — still linear sum of  $x_1, x_2$

$$h_i' = \phi(w_i^T x_i)$$

$y$  is now a non-linear fn. of  $x$ .



Non-linear classification  
OR Regression.



$$(400 \times 200 + 200 \times 10 + 10 \times 1) \rightarrow \text{weights you have to train}.$$

→ backpropagate your loss.  
eg - house price prediction.

$$\left. \begin{array}{l} \vec{h}_1 = \phi(\vec{w}_1^T \vec{x}_i) \\ \vec{h}_2 = \phi(\vec{w}_2^T \vec{h}_1) \\ y = \phi(\vec{w}_3^T \vec{h}_2) \end{array} \right\}$$

$$L = \sum (y_i - \hat{y}_i)^2$$

minimize this loss by changing weights.

→ Use gradient descent.

$$w_{t+1} \rightarrow w_t + \lambda \cdot \frac{\partial L}{\partial w} \quad \leftarrow \text{gradient descent.}$$

[use chain rule].

$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial w} \quad \leftarrow$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w_3}$$

$$= \sum_{\text{over batch}} [2(y - \hat{y}) \cdot \frac{\partial \phi(w_3^T h_2)}{\partial w_3}]$$

## Backpropagation

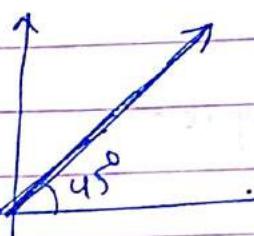
→ gradient descent of with chain rule.



## Neural Netw

→ Gradient start at loss fn & backpropagate.

⇒ Slides.



ReLU Activation fn.

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

piecewise differential.

→ we can manage with these

$\left\{ \begin{array}{l} \text{If p} \rightarrow \text{extract features} \\ \rightarrow \text{learn a classifier.} \end{array} \right.$

SMAI (22/03).

## MLP (Multilayered Perception)

$$x \quad w_1^T x \quad \phi(w_1^T x) \quad w_2^T z \quad \phi(w_2^T z) \quad y \quad \frac{1}{2}(y - \hat{y})$$

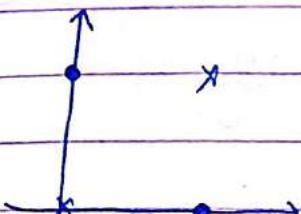
train  $w_1$  &  $w_2$ .

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w_2} = (y - \hat{y}) \cdot \phi'(w_2^T z) \cdot z$$

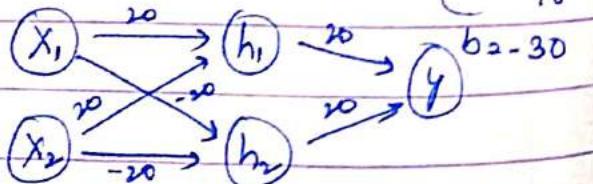
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w_1} = (y - \hat{y}) \cdot \phi'(w_2^T z) \cdot w_2 \cdot \phi'(w_1^T x) \cdot x \cdot w_1$$

XOR

0	0	0
0	1	1
1	0	1
1	1	0



$$\vec{x}_1 = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \vec{h} = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}$$



$$b = 30$$

$$y = \phi(20h_1 + 20h_2 + 30)$$

$$h_1 = \phi(20x_1 + 20x_2 - 10)$$

$$h_2 = \phi(-20x_1 - 20x_2 + 30)$$

$$h_1 \\ \phi(20x_0 + 20x_0 - 10) \approx 0 \\ \phi(0 + 20 - 10) \approx 1$$

$$h_2 \\ \phi(-20x_0 + -20x_0 + 30) \approx 1 \\ \phi(0 - 20 + 30) \approx 1$$

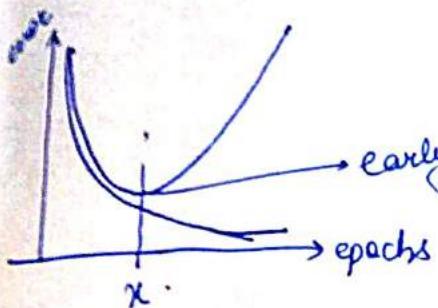
$$y \approx 0.$$

$$\approx 0.1$$

$$y = 20h_1 + 20h_2 - 30.$$

$$y = 200x_1 + 200x_2 + 100 \rightarrow \text{line.}$$

XOR problem can't be solved by a line.



### Loss function

① Regression  $\rightarrow \frac{1}{N} \times \frac{1}{2} \sum_i (y - \hat{y})^2 - \text{MSE} . \quad (\text{MSE} \rightarrow \text{avg.})$

$$\sum_i |y - \hat{y}| - \text{MAE}$$

$$\text{KL divergence} = \sum -p_i \log \frac{q_i}{p_i}$$

$\Rightarrow$  regularization terms can be added.

$$\frac{1}{2} \sum_i (y - \hat{y})^2 + \lambda \sum_{i=1}^N w_i^2$$

### ② Classification

Gross entropy  $\rightarrow$  one hot encoding.

$$\begin{matrix} \text{class 1} & \text{class 2} & \dots \\ \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} & \dots \end{matrix}$$

size of output layer = no. of labels

$$-\log p_i$$

$i = \text{ground truth}$

$$P = \begin{bmatrix} p_0 \\ p_1 \\ \vdots \\ p_{n-1} \end{bmatrix}$$

$$p_0 + p_1 + \dots + p_{n-1} = 1.$$

$\log 0 = -\infty$  ;  $\log 1 = 0$  .  
 ↓  
 max. loss.

$$\text{Binary Cross Entropy} = y \log p_i + (1-y) \log(1-p_i)$$

### SoftMax

$$\begin{bmatrix} q_{10} \\ q_{11} \\ \vdots \\ q_{1m} \end{bmatrix} \xrightarrow{\text{Softmax}} \begin{bmatrix} p_0 \\ p_1 \\ \vdots \\ p_{n-1} \end{bmatrix}$$

$$p_0 = \frac{e^{q_0}}{e^{q_0} + e^{q_1} + \dots + e^{q_{n-1}}} \xrightarrow{\text{softmax}}$$

↳ used in multiclass setting.  
 multiple binary classification  
 multiple class true.

- Binary case doesn't need softmax, we can use sigmoid.
- Sigmoid & softmax used mainly as activation fn.

$$p = \sigma(v)$$

$$p = \text{softmax}(v) = \arg\max_i(p_i^*)$$

0.5
7
5
2
0.3

after softmax. →

0.87
0.118
0.0058

} large diff.

SMAI (25/03/21)

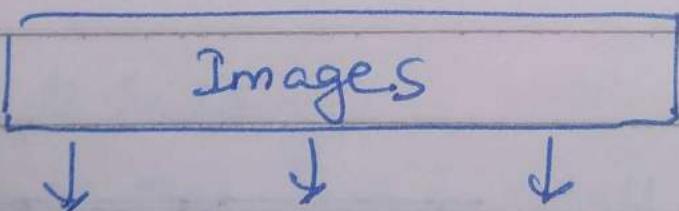
## CNN

Image  $\rightarrow$  3-D array (H-w-L R-G-B).

Image  $\rightarrow$  Flatten it  $\downarrow$  use SVM, etc.  
large size.

Other features  $\rightarrow$   $\begin{bmatrix} \text{Min Red.} \\ \text{Max Red} \\ \text{Mean Red.} \\ | \\ ; \end{bmatrix}$   $\rightarrow$  Same class  
dist. can be  
diff. Some; diff class  
some dist.

## Higher level of Representations: Bag of Words



Take patches.



draw histograms (of fixed dim.).

CNN  $\rightarrow$  ImageNet

## Dense Connections

$$w^T x \quad [ ]$$

$$w[ ] \left\{ \begin{array}{c} \diagup \\ \diagdown \end{array} \right\} \quad x \rightarrow [ ]$$

→ densely connected.

## Convolutional Connection

PPT.

1	0	2
---	---	---

filter ( $1 \times 3$ )

4	1	7	3	5
---	---	---	---	---

$I/P (1 \times 5)$



18	7	17
----	---	----

$O/P (1 \times 3)$

→ padding may or may not be done.

e.g.

1	2	0	1
0	4	7	2
3	1	4	0
1	2	0	2

$4 \times 4$

2	1	2
7	0	0
4	1	3

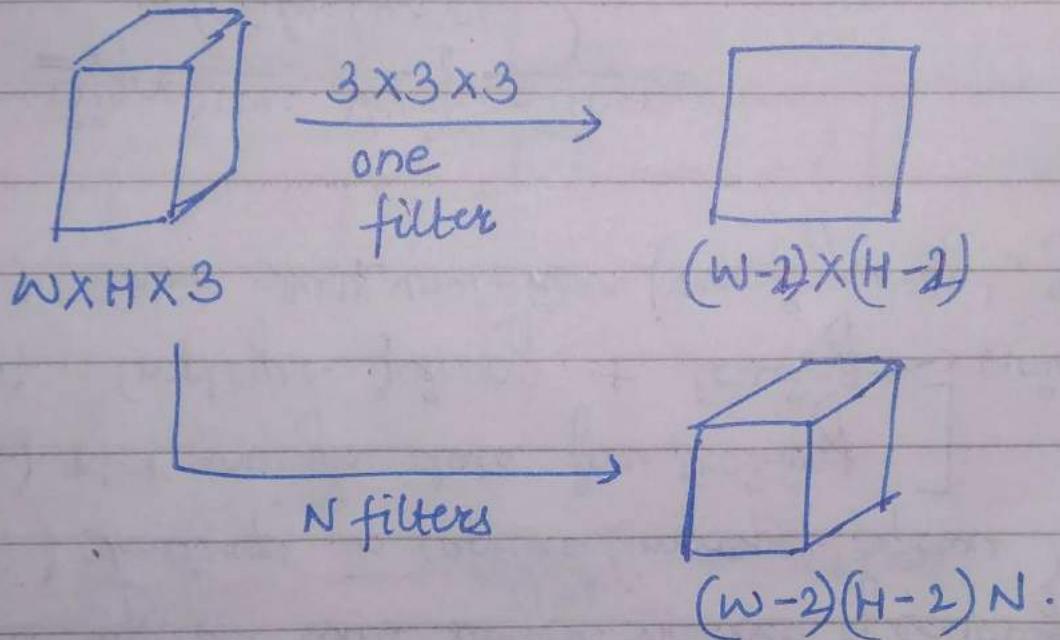
$\xrightarrow{3 \times 3}$

29	42
----	----

$2 \times 2$

→ taking  $3 \times 3$  filter & traversing it over image.

$\rightarrow N$  filter of  $3 \times 3$   $\xrightarrow[\text{padding}]{\text{no}} 2 \times 2 \times N$  o/p  $\downarrow$  depth  
 $\xrightarrow{\text{padding}} 4 \times 4 \times N$  o/p.



idea of convolution layer.

$\rightarrow$  Stride & Padding.

$\rightarrow$  Pooling — take max. in small window.

# max-pooling ; min-pooling ; arg-pooling

Sia

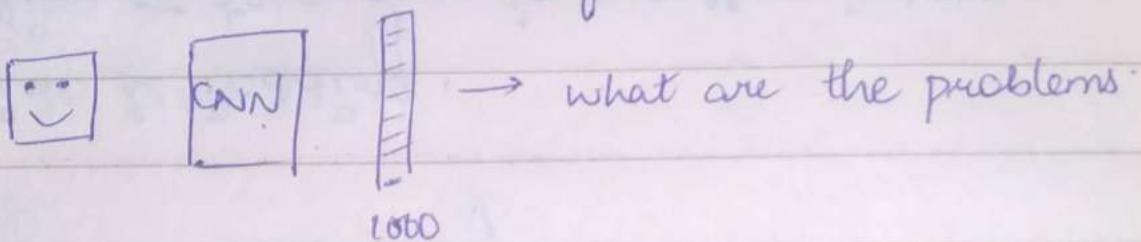
SMAI (01/04)

## Siamese Networks + Metric learning

→ Face recognition based entry system for company.

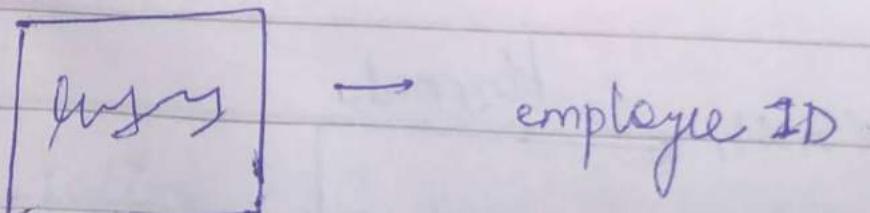
[Suppose 1000 employees]

Sol 1 — 1000 class classification



- (i) difficult with larger number of employees
- (ii) challenges in capture (improper capture) → agnostic
- (iii) Not enough data for training
- (iv) Employee joining + leaving → train again every time
- (v) Intra class variation (removing beard, goggle etc.)

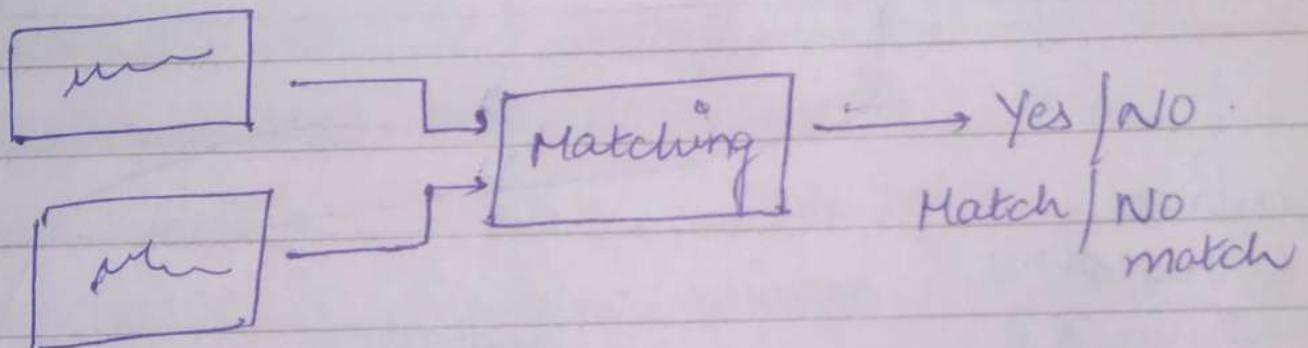
→ Signature Verification  
(- in france)



15 million classification prob.

↓  
not doable

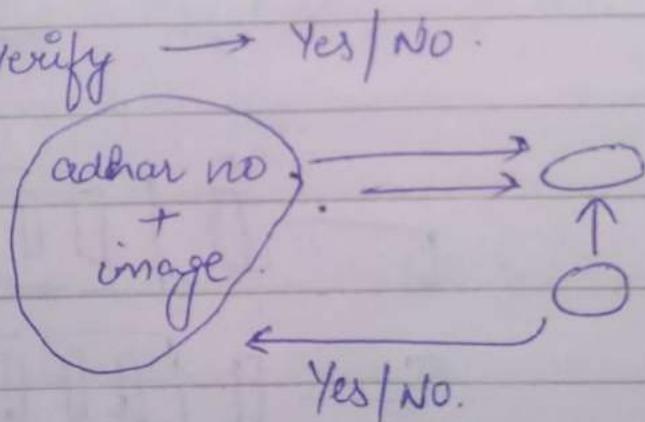
→ problem is ~~not~~ now verification problem.



→ Aadhar

during verification  
only thumbs

VIDAI



During registration

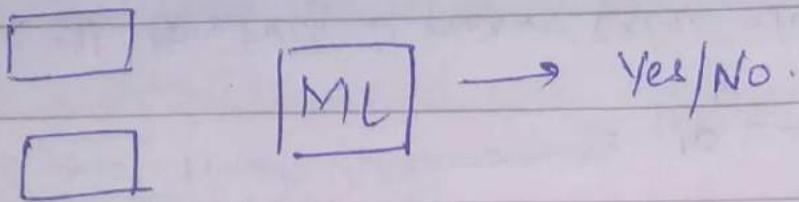
→ 2 thumbs, 8 fingers, IRIS - ? why all this data

Thumb recognition - 99.99999 %. accurate

(1 thumb →  $10^5$  - 1 mistake) →  $10^5$  - 1 mistake.

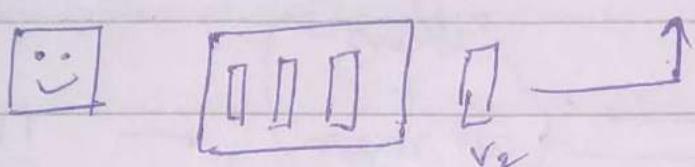
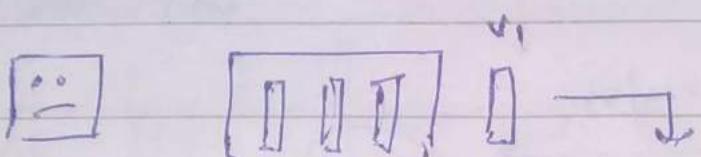
[ take all to check if person exist already in db ]  
while registering decreasing prob. of error.

## Verification



Hp

## Face verification



$y=1 \rightarrow$  faces are of same

identity  $d(v_1, v_2)$  → small  
(dist)

$y=0 \rightarrow$  faces are of diff.  
identity

$d(v_1, v_2)$  → large

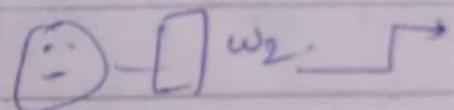
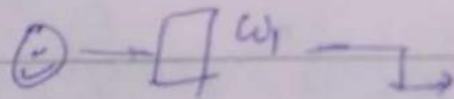
&  
more than a  
margin.

$$\text{loss fn.} = \min \left( y \cdot d^2 + (1-y) \max(0, m-d) \right)$$

→ contrasting loss fn.

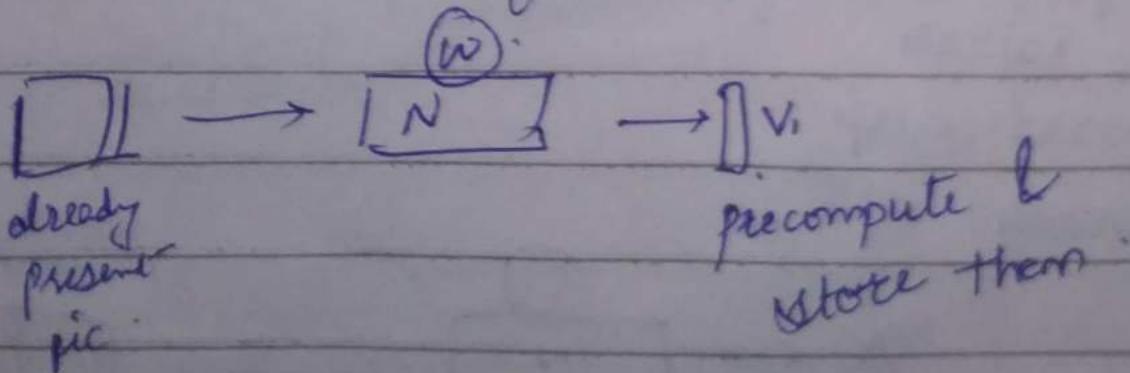
- idea of contrastive learning — metric learning
- post processing had to be done earlier
- Now deep learning used — metric learning
- now we are transforming a face into a vector in that vector space (metric), faces of same identity are closer to each other  
 (irrespective of intra class variation) — that's why metric learning

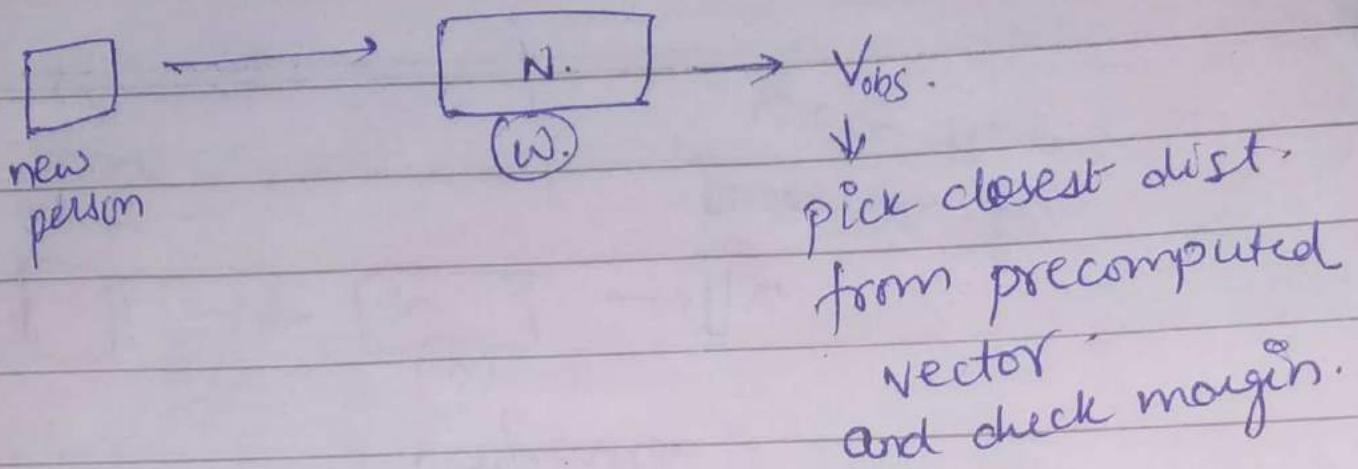
→ Shared network



update  $w_1$  &  $w_2 \rightarrow w = w_1 + w_2 / 2$  i.e mean.

Inference : 1000 employee problem  
 (pre-trained network for face verification)





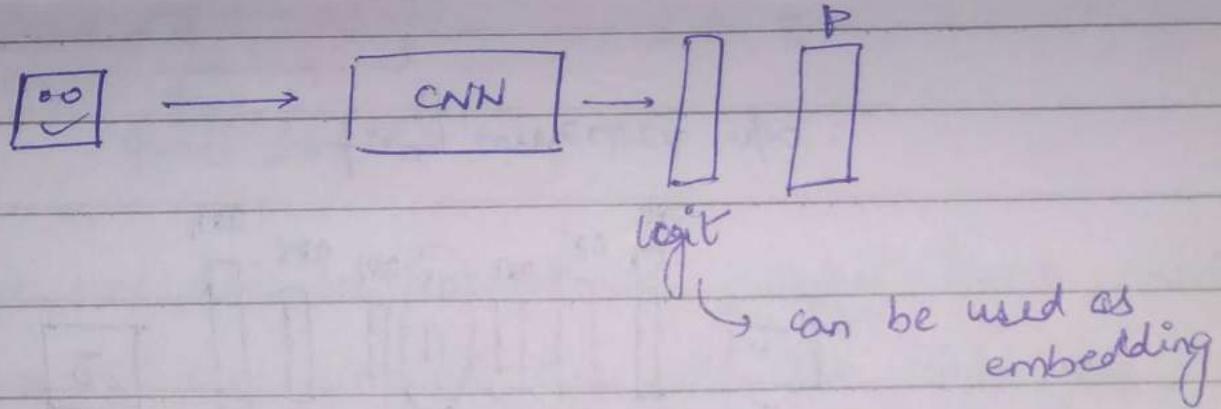
→ One forward pass + dist computation only.

req. req.

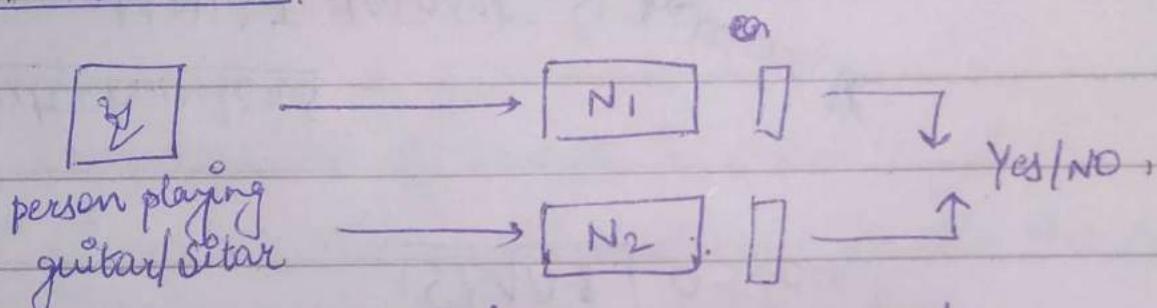
Hard Negative Mining - use negatives which are harder  
 to distinguish, training will be much  
 stronger.

→ Celebrity dataset not work correctly on  
 real world data due to biases

- clarity of photo, smiling, makeup on  
 These biases can drop accuracy on real  
 world data.



## Multi Model.

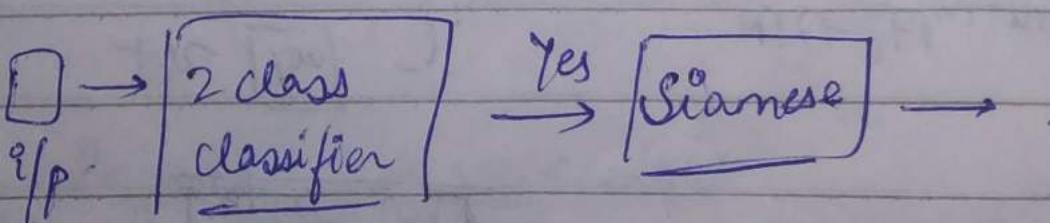


If image & caption match — Yes  
else — No.

Q. What if instead of human, photo is placed?

~~instead of Ans ①~~ → depth sensor reading (depth measurement)

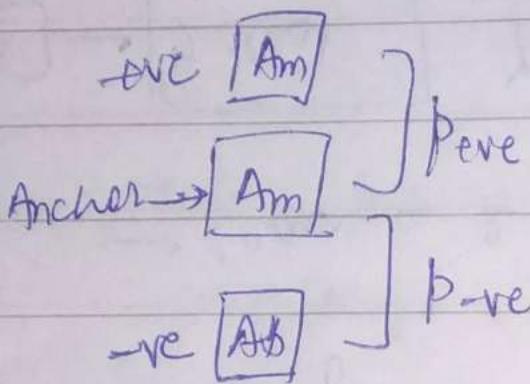
→ Learn a 2 class classifier (real person vs. photograph.)



## → Person Re-identification

Matching patches - Shift (best algo)

Triplet loss



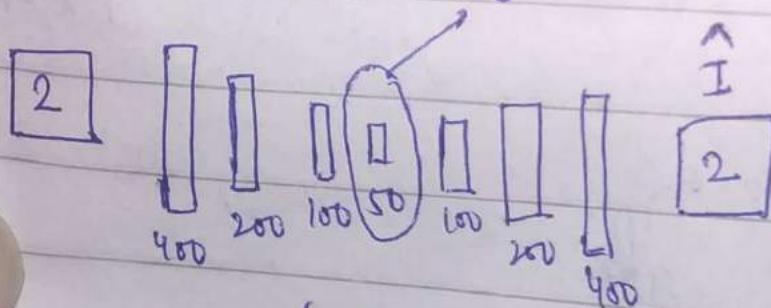
dist. b/w non-matching pair  
should be atleast a  
margin more than matching  
pair.

(Contrastive & triplet gives almost similar accuraccy)

## SMAI (05/04)

### Autoencoders

MNIST dataset bottleneck



fully connected n/w.

$$\text{loss} = \mathbb{E}(\hat{\mathbf{x}} - \mathbf{x})$$

$$\hat{\mathbf{x}} \approx \mathbf{x}$$

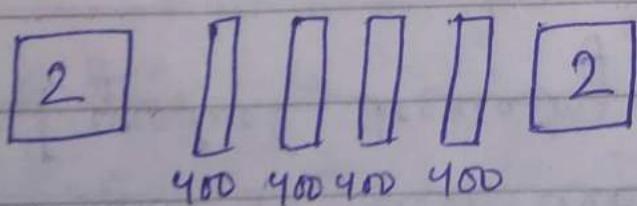
quite small

→ if we are able to reconstruct image back, that means 50 dim. vector has all the info req. to construct the image.

- Compression of image is done.
- Smaller the bottleneck — greater the compression.
- Instead of 900 dim. image — apply ML on 50 dim.
- PCA — linear

Autoencoder — non-linear (activation fn. can be non-linear).

### Network



→ just passing of info.  
→ Net should be able to do it.

← → identity mapping

### Benefits of Autoencoders

- compression
- dim. reduction — feature extraction  
(good feature extractor)

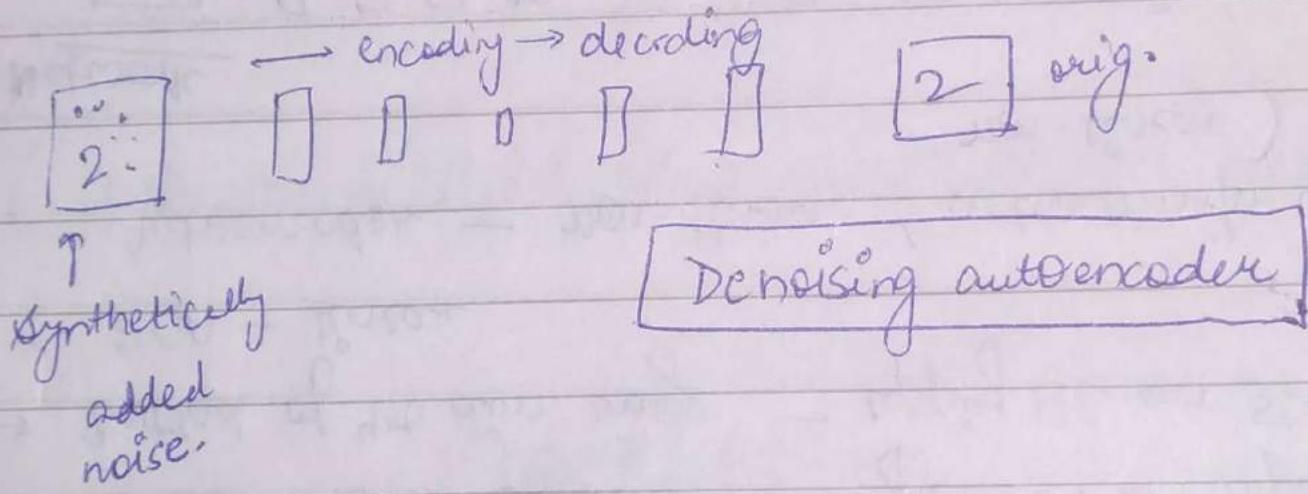
- PCA not always good feature extractor
- Retrieval done using kNN mostly.  
(if less dim. retrieval very fast)

→  $784 \rightarrow 30 \rightarrow 784$

$$784 \rightarrow 1000 \rightarrow 500 \rightarrow 250 \rightarrow 30 \rightarrow 250 \rightarrow 500 \rightarrow 1000 \rightarrow 784$$

2<sup>nd</sup> case i.e. deep encoder is better.

[gradual decreasing and increasing better.]



→ Self Supervision : people create dataset

grey  $\rightarrow$  color.

$\downarrow$   
(weighted sum of color gives grey)

[creating dataset easy.]

→ supervised / self-supervised . autoencoders .

## GANs (Generative Modeling)

→ Given training data , generate new samples from same distribution .

### GANs (Application of Game Theory)

→ generative adversarial networks

→ A game b/w two players -

discriminator D

Generator G.

→ D tries to discriminate b/w

a sample from data dist .

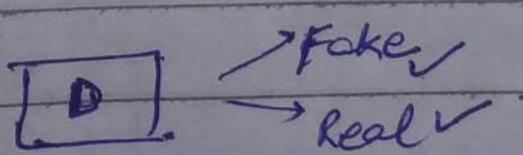
a sample from generator G .

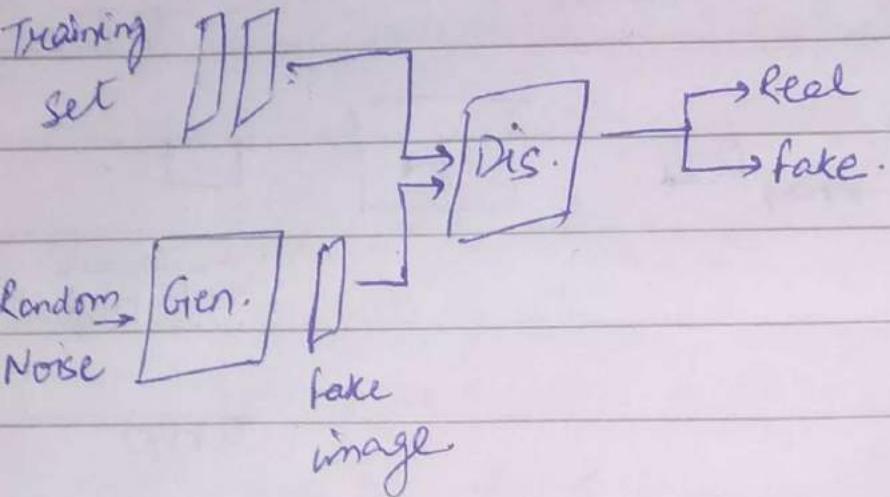
→ G tries to "trick" D by generating samples that are hard for D to distinguish from data .

### FNG



→ Imp .





→ no explicit supervision in GAN's  
 freeze gen. → train discriminator → freeze dis.  
 → train gen.

[continue this way]

# Key Idea - Slides.

→ GAN difficult to train.

SMAI (12/04).

Word2Vec → 1) lower dimensional.

Mango =  $\begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}_{300}$  2) semantics.

Hotel      Motel  
 $\begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}, \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}$

Hotel X Motel  $\approx 20$  (noise)

how to create embedding?

→ curate dataset - fill in the blanks.

2 layer MLP (paper)



dictionary  
size=dok

2 layers  
10k.      O/P(ioh)

Softmax → cross entropy

(one hot  
embedding)

$$W = 10,000 \times 300$$

$$U = 300 \times 10,000$$

$$z = w^T x \quad z \rightarrow \text{hidden layer}$$

$$y = u^T z \quad \downarrow \text{vector}$$

after training it  
becomes wordvec.

→ Skip gram model. (slightly beneficial if word occurs less)

→ CBOW Model

King - man + woman  $\rightarrow$  queen.

[ ] [ ] [ ] [ ]

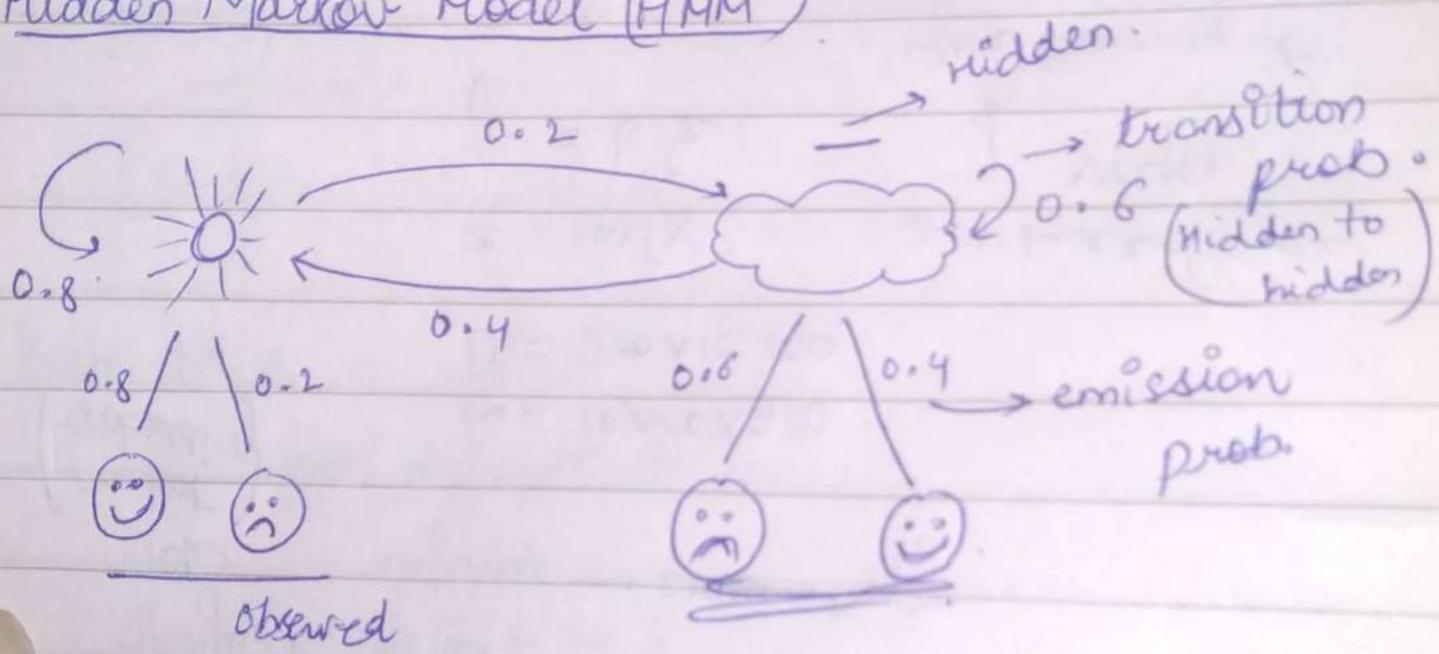
$$S_1 \approx S_2 \approx S_3$$

$$\text{India} - \text{delhi} (S_1) \quad \text{france} - \text{paris} (S_2)$$
$$\text{nepal} - \text{kathm} (S_3) \quad [ ] - [ ] = [ ] \quad [ ] - [ ] = [ ] \quad [ ] - [ ] = [ ]$$

Glove, Word2Vec, Bert, Elmo, USE } for word  
 ↓  
 google embedding.

{ techniques used are - supervised -  
 data comes easy / labels comes free - self supervised }

## Hidden Markov Model (HMM)



→ Given seq. of observations, I want to predict seq. of hidden variables.

→ HMM — predicting what you are not observing

How do I get these prob? — emission & transition

eg- SSSS RRR SSSS RR SS  
G H H H G G H G H H H  
Trans. prob.

$$S \rightarrow S = 8 = 0.8 \quad R \rightarrow R = 3 = 0.6$$
$$R \rightarrow S = 2 = 0.2 \quad S \rightarrow R = 2 = 0.4$$

emm. prob.

$$S \rightarrow H = P(H|S) = 0.8$$

$$S \rightarrow G = P(G|S) = 0.2$$

$$P(H|R) = 0.4$$

$$P(G|R) = 0.6$$

## Viterbi Algorithm

	M	T	W	Th	F	Sat	most likely config. to reach this box.
Y <sub>3</sub>	H	H	G	G	G	H.	
Z <sub>4</sub>	$\frac{2}{3} \times 0.8 = 0.533$	$0.533 \times 0.8 = 0.4264$	$0.0546$	$0.0087$	$0.0014$	$0.0017$	
Y <sub>4</sub>	$\frac{1}{3} \times 0.4 = 0.133$	$0.533 \times 0.2 = 0.1066$	$0.041$	$0.0147$	$0.0053$	$0.0013$	$= 0.341$

Bunny on Tuesday  $\rightarrow \max. \left\{ \begin{array}{l} 0.533 \times 0.8 \times 0.8, \\ 0.133 \times 0.4 \times 0.4 \end{array} \right.$

Rainy on Tuesday  $\rightarrow \max. \left\{ \begin{array}{l} 0.533 \times 0.2 \times 0.4 \\ 0.133 \times 0.6 \times 0.4 \end{array} \right. - 0.043$

$$\text{Sunny on } w = \max \cdot \begin{cases} 0.341 \times 0.8 \times 0.2 \\ 0.043 \times 0.2 \times 0.4 \end{cases}$$

$$R \text{ on } w = \max \cdot \begin{cases} 0.341 \times 0.2 \times 0.6 \\ 0.043 \times 0.6 \times 0.6 \end{cases}$$

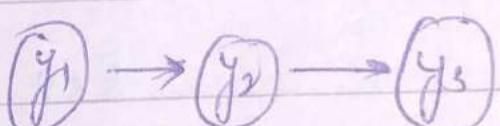
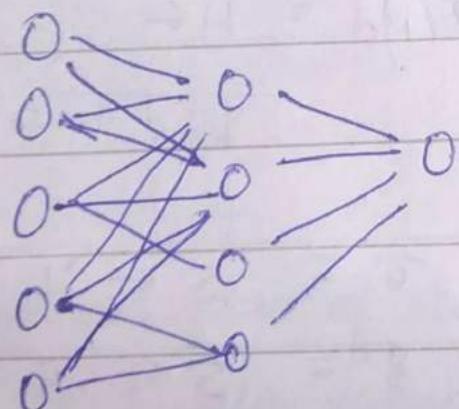
O/P prediction — SSSRRRS

HMM → object tracking  
 → automatic cinematographic  
 (camera selection)

SMA1 (15/04)

Recurrent Neural Networks (RNN's)

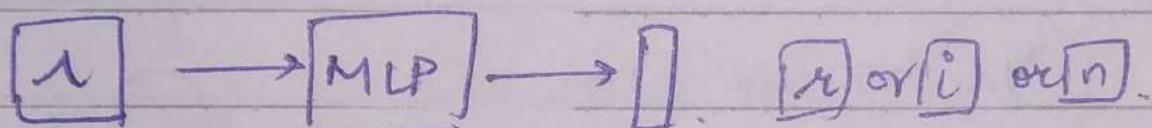
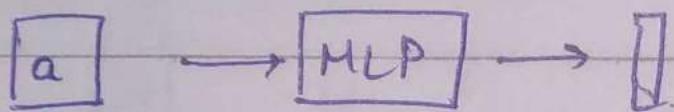
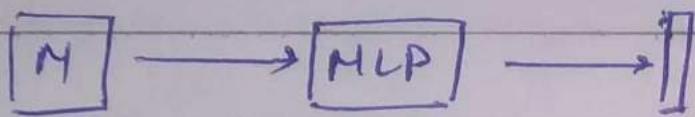
MLP's



$$\begin{aligned} y_2 &= \phi(w_1^T y_1) \\ y_3 &= \phi(w_2^T y_2) \end{aligned} \quad \left. \right\}$$

one sample at a time

## Machine



- $\rightarrow$  one sample at a time       $\left\{ \begin{array}{l} \text{prob. is} \\ \text{MLP} \end{array} \right.$   
 $\rightarrow$  fixed dimension

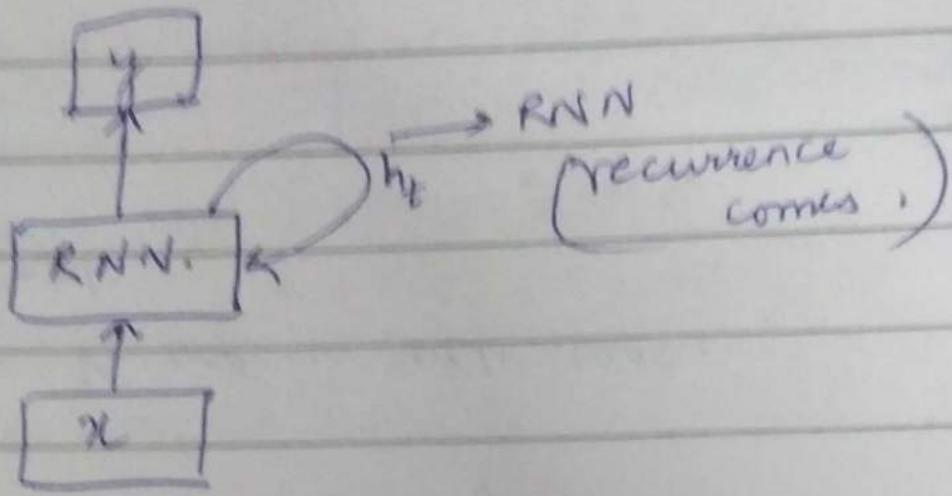
## Sentiment Analysis

$\left\{ \begin{array}{l} \rightarrow \text{It was a good movie.} \\ \rightarrow \text{B - - - -} \\ \qquad \qquad \qquad \text{--- but I did not like the movie} \end{array} \right.$

- $\rightarrow$  The food was good, not bad at all.       $\left\{ \begin{array}{l} \text{Bag of} \\ \text{words} \end{array} \right.$   
 $\rightarrow$  The food was bad, not good at all.       $\left\{ \begin{array}{l} \text{similar} \end{array} \right.$

Order is imp. in such case.

So, Need RNN



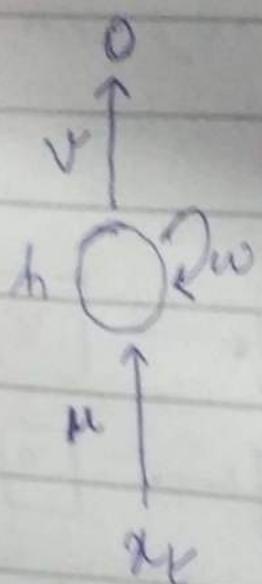
ML used RNNs.

→ seq. to seq. predictions:

(long. translations, unimodal or multimodal us, ASR)

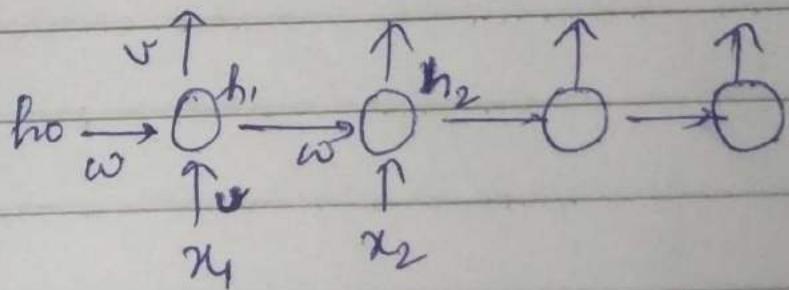
→ time series prediction

(→ stock price prediction etc.  
prediction based on prev. data.)



$$h_t = f(Ux_t + Wh_{t-1})$$

$$o_t = \text{softmax}(Vh_t)$$



$$h_1 = f(u_{x_1} + w h_0)$$

$$o_1 = \text{softmax}(v h_1)$$

$$h_2 = f(u_{x_2} + w h_1)$$

$$o_2 = \text{softmax}(v h_2)$$

} current o/p will be based on current i/p &  
 } past observation.

$$\boxed{\frac{\partial L}{\partial w} = \sum_t \frac{\partial L}{\partial o_t} \frac{\partial o_t}{\partial w}}$$

→ RNN stores short term info.

LSTM's

Stacked RNN