



(रज. नं. क-१९८-राजकोट ता. १७-प-२००१)

# શ્રી ઉમિયા યુવા ચેરીટેબલ ટ્રસ્ટ-રાજકોટ

૧૫૨ - બેકબોન શોપીંગ સેન્ટર, માયાએની ચોક,  
ચંદ્રેશાનગર મેઈન રોડ, રાજકોટ-૪. ફોન : ૨૩૬૭૧૩૬



hucenmall@vsnl.com

## S P E C I A L

**HI-BOND** <sup>TM</sup>  
c e m e n t

Bonding Nation...

WEBSITE :- [www.hibondcement.com](http://www.hibondcement.com)

EMAIL :- [sales@hibondcement.com](mailto:sales@hibondcement.com)

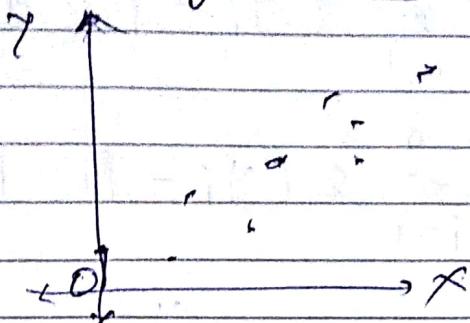
I N D E X

Name : Kaushal P. Shah Subject : Machine Learning  
Std : \_\_\_\_\_ Div. : \_\_\_\_\_ Roll No. : \_\_\_\_\_  
School / College : \_\_\_\_\_

# MACHINE LEARNING

Page No.		
Date		

\* Linear Regression (Line fitting)



$h(x_i) = \beta_0 + \beta_1 x_i$ , where  $h(x_i)$  is predicted value.

$$\text{so error } (\epsilon_i) = \underbrace{(y_i)}_{\text{original value}} - h(x_i)$$

→ Our aim is to minimize this error, ( $\epsilon_i$ )

→ So mean squared error ( $J(\beta_0, \beta_1)$ ) is

$$J(\beta_0, \beta_1) = \frac{1}{2n} \sum_{i=1}^n \epsilon_i^2$$

$$= \frac{1}{2n} \sum_{i=1}^n (y_i - h(x_i))^2$$

→ To find Global minima (bcz we are trying to minimize error  $\epsilon_i$ ), we need to differentiate above eqn.

$$\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_0} = \frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1} = 0$$

$$J(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial J(\beta_0, \beta_1)}{\partial (\beta_0)} = \frac{1}{n} \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-1) \rightarrow ①$$

and

~~$$\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1} = -\frac{2}{n} \sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i)$$~~

$$\rightarrow ②$$

→ Comparing both to zero and solving both eq's.

$$\rightarrow -\frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$-\frac{2}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\rightarrow \therefore \sum_{i=1}^n y_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0 \rightarrow ①$$

and  $\sum_{i=1}^n y_i x_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0 \rightarrow ②$

$$\rightarrow \therefore \sum_{i=1}^n y_i - \beta_0 n - \beta_1 \sum_{i=1}^n x_i = 0$$

$$\sum_{i=1}^n y_i x_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0$$

$$\sum y \sum x - n \beta_0 \sum x - \beta_1 (\sum x)^2 = 0$$

$$\begin{matrix} n \sum xy & - n \beta_0 \sum x & - n \beta_1 \sum x^2 \\ - & + & + \end{matrix} = 0$$

$$\sum x \sum y - n \sum xy - \beta_1 (\sum x)^2 + n \beta_1 \sum x^2 = 0$$

$$\sum x \sum y - n \sum xy + \beta_1 (n \sum x^2 - (\sum x)^2) = 0$$

$$\beta_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$\rightarrow$  writing  $\frac{\sum a}{n} = \bar{a}$

$$\beta_1 = \frac{\cancel{\sum xy}}{n^2} - \frac{\sum x \sum y}{n^2}$$

$$\frac{\sum x^2 - (\sum x)^2}{n}$$

$$\boxed{\beta_1 = \frac{\sum xy - n \bar{x} \bar{y}}{\sum x^2 - n (\bar{x})^2}}$$

$$\left[ \frac{\sum x^2 - n \bar{x}^2}{n} \right] \& \boxed{\beta_0 = \bar{y} - \beta_1 \bar{x}}$$

→ Similar for multiple features (Multiple variables)

~~Ex. 1~~

$$J(\beta_0, \beta_1, \beta_2, \dots, \beta_m) = \frac{1}{2n} \sum_{i=1}^n (h(x_i) - y_i)^2$$

- We can similarly derive  $\beta_0, \beta_1, \dots, \beta_m$  as done previously.
- But such formula are not used in practice.

\* Gradient Descent with one variable cost

- Have some function  $J(\beta_0, \beta_1)$  or  $J(\theta_0, \theta_1)$  & we want to minimize it.
- Start with some initial  $\theta_0, \theta_1$
- Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we get minimum

\* Algo.

repeat until convergence {

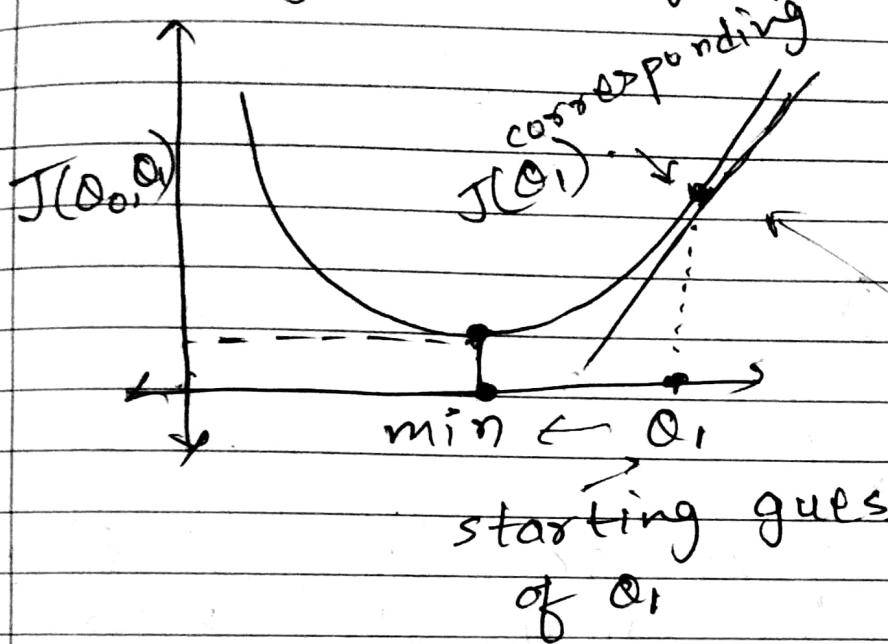
$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

}

where,  $\alpha$  is learning rate or step size.

→ How this works?

→ Assume  $\theta_0 = 0$  for simplicity (for getting 2-D graph of  $J(\theta_0, \theta_1)$ ).



$$\rightarrow \text{now, } \theta_1 = \theta_1 - \alpha \underbrace{\frac{d}{d\theta_1} J(\theta_1)}_{\text{slope at } \theta_1}$$

$\rightarrow$  as seen from figure, slope is +ve.

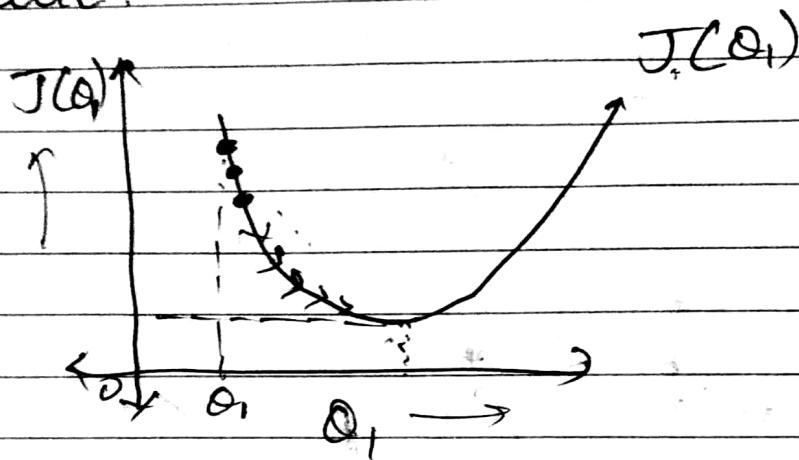
so,

$$\theta_1 = \theta_1 - \alpha (+\text{ve value})$$

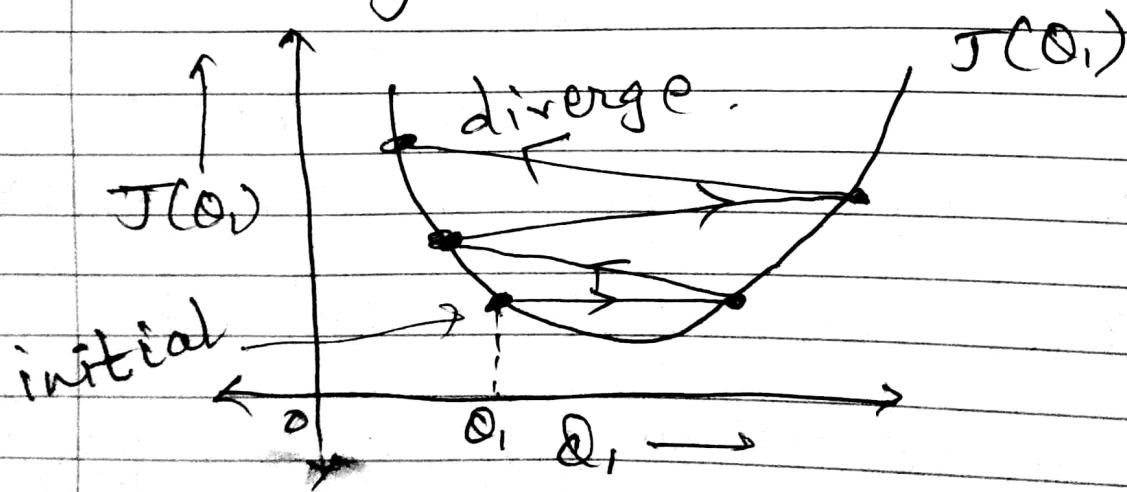
so  $\theta_1$  is decreasing, or going towards the minimum value

$\rightarrow$   $\theta_1$  will increase if we had ~~guessed~~ initialized it at LHS of min. value because slope will be -ve there & from above eq',  $\theta_1$  will increase or go to min. value.

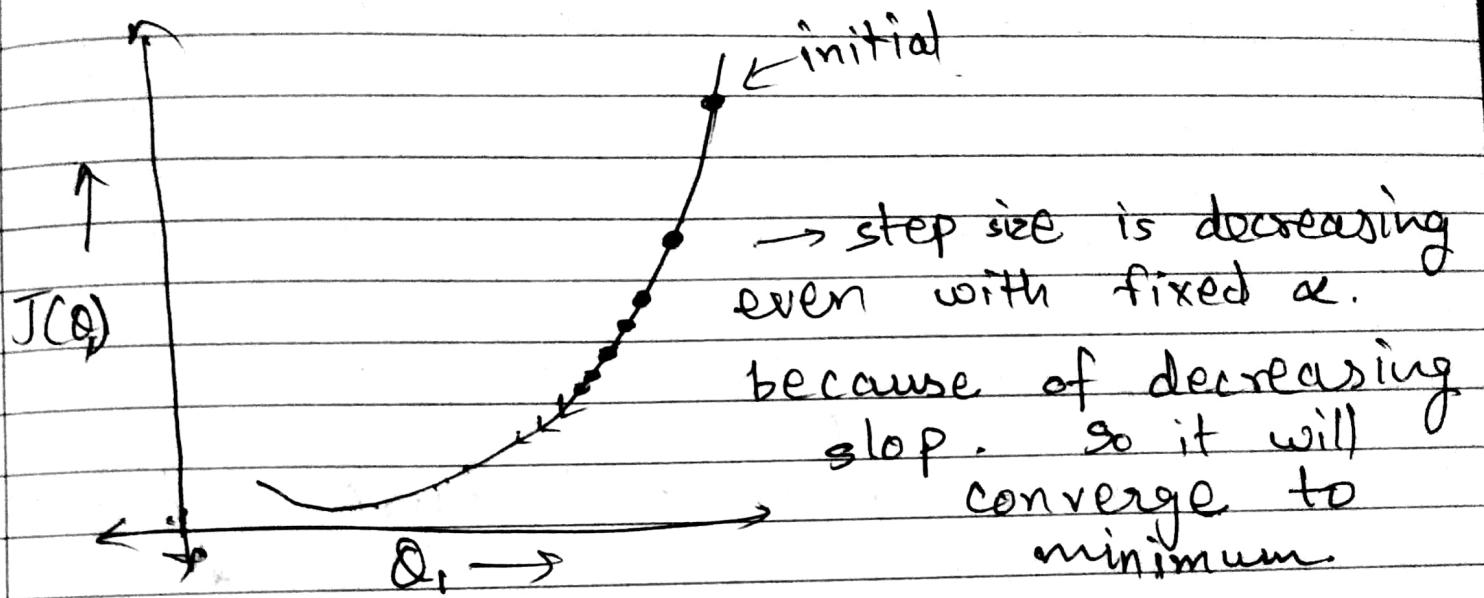
→ If  $\alpha$  is too small, steps will be very small, so more iterations to reach minimum value.



↔ And if  $\alpha$  is too large, it may overshoot the min. It may fail to converge or even diverge.



→ Gradient descent can converge to local minima even with fixed  $\alpha$ .



→ "Batch" Gradient Descent + Each step of gradient descent uses all the ~~the~~ training examples.

## Lecture-2

Page No.	
Date	

### \* Linear Regression.

$x_1$  = size of house

$x_2$  = # of bedrooms, then

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

if we define  $x_0 = 1$ , then.

$$h_{\theta}(x) = \sum_{i=0}^2 \theta_i x_i = \theta^T x.$$

if  $n$  is # of features, (here  $n=2$ ), then

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x.$$

→  $\theta$ 's are called # of parameters, which we are going to learn,

→ now error  $E = (h_{\theta}(x) - y)$

where  $h_{\theta}(x)$  = predicted output or hypothesis  
and  $y$  is actual output.

→ Mean squared error can be given as

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

→ our goal is  $\min_{\theta} J(\theta)$  (minimize  $J(\theta)$  over param  $\theta$ )

- Batch Gradient descent:- Update  $\theta_i$  after iterating whole training set.
- Stochastic GD:- Update ~~all~~ all  $\theta_i$  after for each training examples so this is more faster than Batch Gradient Descent.
- now, writing all terms in matrix (vector form),

$$\nabla_{\theta} J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix} \in \mathbb{R}^{n+1}$$

- so, Batch GD → iterations can be given by

$$\theta = \theta - \alpha \nabla_{\theta} J$$

vector  
 $R^{n+1}$

$(x^{(i)})^T$  is nothing but  
a row

Page No.	
Date	

$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \text{ size } n \times m$$

$X$  is matrix of training ex.

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}_{m \times 1}, \vec{\theta}$$
 is vector of params to be learned.

so,  $X\vec{\theta} = \begin{bmatrix} (x^{(1)})^T \vec{\theta} \\ (x^{(2)})^T \vec{\theta} \\ \vdots \\ (x^{(m)})^T \vec{\theta} \end{bmatrix} = \begin{bmatrix} h_{\vec{\theta}}(x^{(1)}) \\ \vdots \\ h_{\vec{\theta}}(x^{(m)}) \end{bmatrix}$

this entire row of  $X$ .

$$\rightarrow \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}_{m \times 1}, \text{ where } \vec{y} \text{ is vector of actual outputs.}$$

$$X\vec{\theta} - \vec{y} = \begin{bmatrix} h_{\vec{\theta}}(x^{(1)}) - y^{(1)} \\ \vdots \\ h_{\vec{\theta}}(x^{(m)}) - y^{(m)} \end{bmatrix}_{m \times 1}$$

vector of  $m \times 1$

$\rightarrow$  If  $\vec{z}$  is any vector,

$$z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}, \text{ then}$$

$$z^T z = \sum_{i=1}^m [z_1 \ z_2 \ \dots \ z_m] \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$$

$$Z^T Z = \sum_{i=1}^m Z_i^2 \quad \text{for any vector } \vec{z}$$

$$\text{so, } \frac{1}{2} (\vec{x}_0 - \vec{y})^T (\vec{x}_0 - \vec{y}) = \frac{1}{2} \sum_{i=1}^m (h(\vec{x}^i) - y^{(i)})^2$$

$$so, J(\theta) = \frac{1}{2} \sum (x\theta - y)^T (x\theta - y)$$

→ Now to minimize  $J(\theta)$ , w.r.t  $\theta$ ,  
 we need to take derivative &  
 compare it to zero.

$$\nabla_{\theta} J(\theta) = \vec{0}$$

$$\nabla_{\theta} \frac{1}{2} (x_0 - y)^T (x_0 - y) = 0$$

$$\therefore \frac{1}{2} \nabla_{\theta} (x_0 - y)^T (x_0 - y) = 0$$

$$\nabla_{x_0} \left( (x_0)^T x_0 - (x_0)^T y - y^T (x_0) + y^T y \right) = 0$$

$$\therefore \nabla_Q (Q^T X Q - Q^T X^T y - y^T X Q - y^T y) = 0$$

→ now  $\sqrt{a}$  is a real number,  
 & trace or tr(a) of real no. is  
 just a real number, so

Page No. \_\_\_\_\_  
Date \_\_\_\_\_

becz it's transpose & tr number  
 bcz taken a real number  
 $\text{tr } A = \text{tr } A^T$

$$\therefore \nabla_Q \text{tr} (\underbrace{Q^T X^T X}_B - \underbrace{Q^T X^T Y}_C - Y^T X Q + Y^T Y) = 0$$

$$\therefore \nabla_Q [\text{tr } Q Q^T X^T X - \text{tr } (Y^T X Q) - \text{tr } (Y^T Y)]$$

now,  $\nabla_Q \text{tr } Q Q^T X^T X = \nabla_Q \text{tr } \underbrace{Q}_A \underbrace{I}_B \underbrace{Q^T X^T X}_C$

$$= \underbrace{X^T X Q}_C \underbrace{I}_B + \underbrace{X^T X Q^T I}_C \underbrace{I}_B$$

also,  $\nabla_Q \text{tr} (\underbrace{Y^T X Q}_B) = \underbrace{X^T Y}_C$

so,  
 $\nabla_Q J(Q) = 0 = [X^T X Q + X^T X Q - X^T Y - X^T Y]$

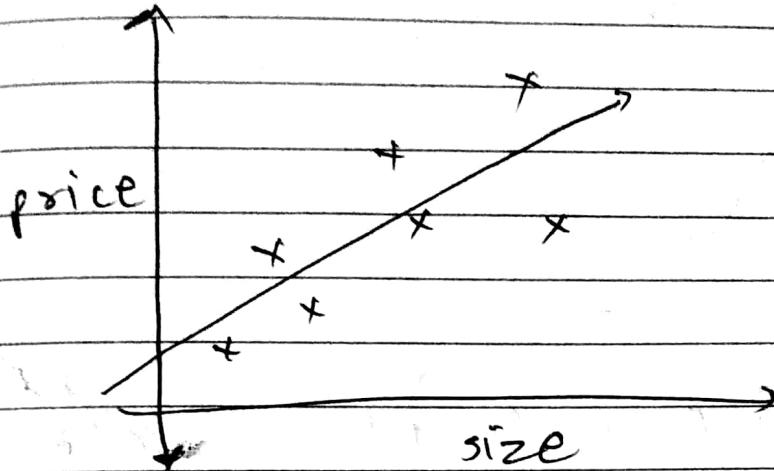
$\therefore X^T X Q = X^T Y \rightarrow \text{Normal eqn}$   
 & solve it for Q

$$\therefore \boxed{Q = (X^T X)^{-1} X^T Y}$$

## Lecture-3

Page No.	
Date	

→ If there is only one feature, for ex. for predicting ~~size~~ prices of houses, if there is only one feature 'size', then

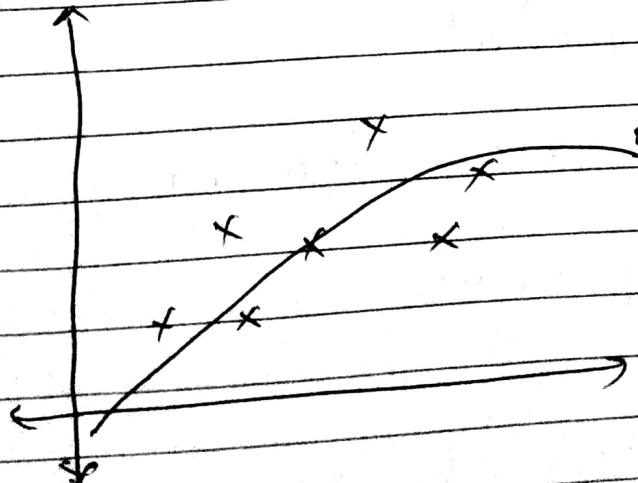


$$\theta_0 + \theta_1 x_1$$

(Underfitting)

→ but if there are two features, size &  $(size)^2$ , then our  $h_0(x)$  would be  $\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$

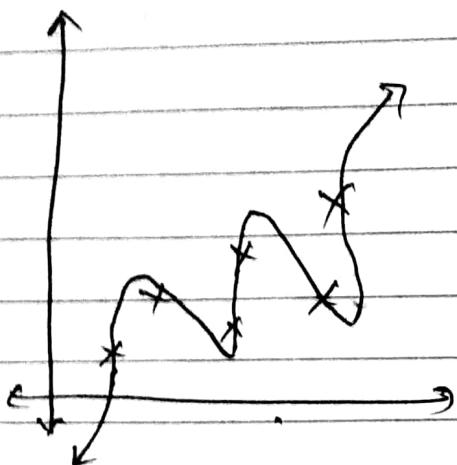
so,



quadratic.

so best fits  
to ~~the~~ model.

→ but if we take too many features, then curve will overfit to data.



$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \dots + \theta_6 x_1^6$$

→ So we need to choose features very carefully.

\* Two types of learning algorithms:

1) "Parametric" learning algo.

- It's defined as an algo. which has fixed no. of parameters, to be fitted to data.

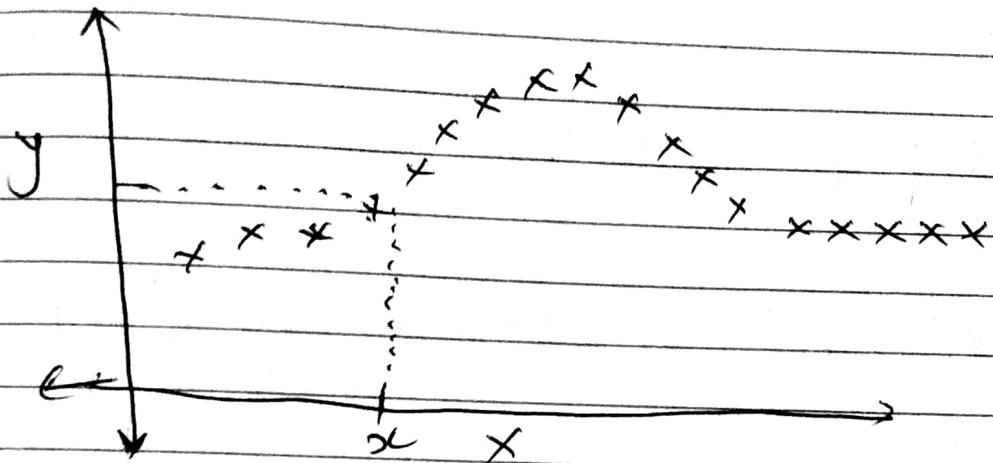
2) "Non parametric" learning algo.

- # of param grows with m  
(size of training set).

\* Locally weighted Regression is one example of "non-param." Algorithm.

(LWR)

\* Locally Weighted Regression / Loess / Lowess.

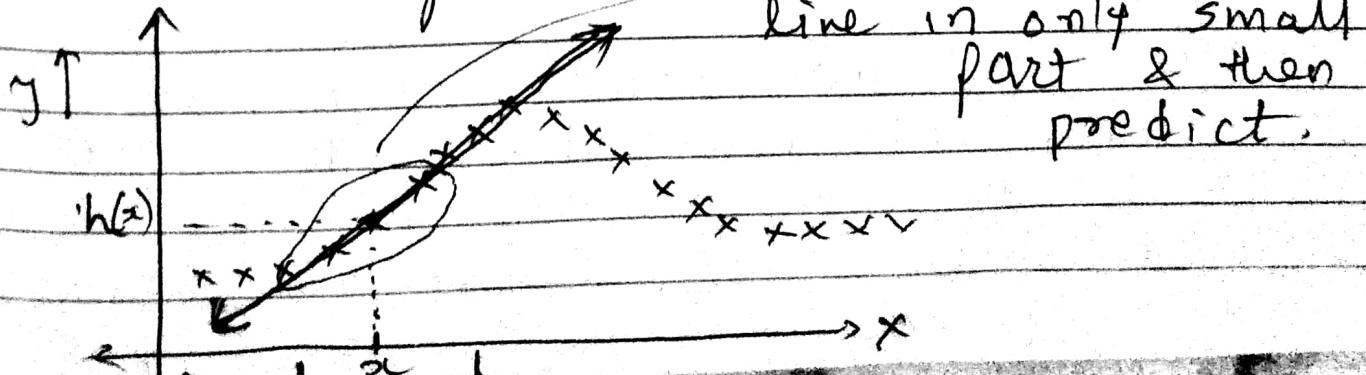


- In above data linear or quadratic function can't fit perfectly.
- To evaluate  $h$  at a certain  $x$ ,

In Linear Regression we do following

- fit  $\theta$  to minimize  $\sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$   
& then return  $\theta^T x$ .

- But in LWR, we apply Linear Regression locally & fit a line locally.



\* LWP: fit  $\theta$  to minimize

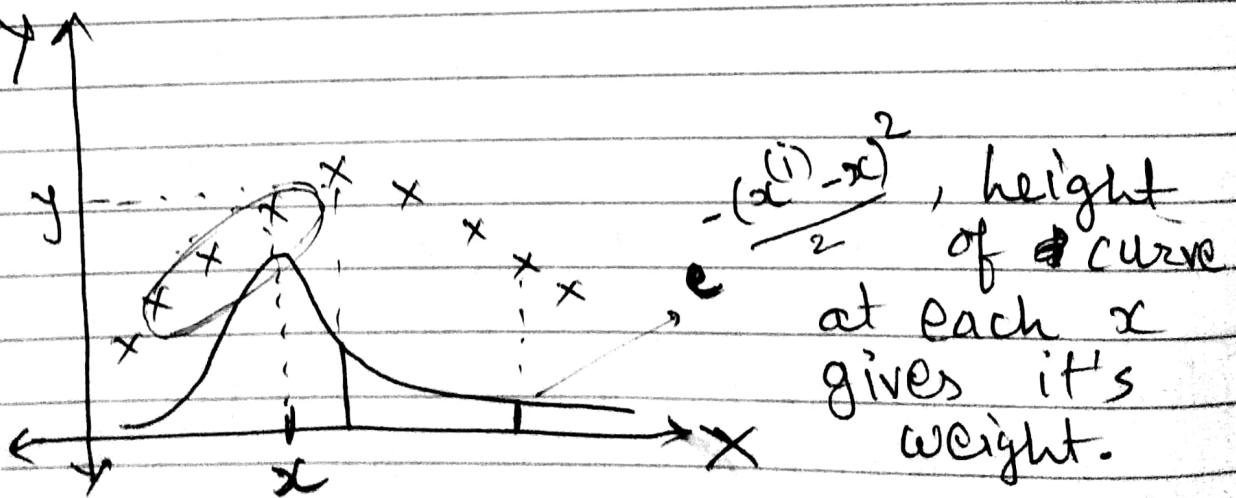
$$\sum_i \omega^{(i)} (y^{(i)} - h_{\theta}(x^{(i)}) )^2$$

$$\omega^{(i)} = \exp \left( - \frac{(x^{(i)} - x)^2}{2} \right).$$

If  $|x^{(i)} - x|$  small, then  $e^0 = 1$   
 $\therefore \omega^{(i)} = 1$   
 $x^{(i)}$  is close to  $x$ .

→ And if  $|x^{(i)} - x|$  is large,  
then  $e^{-\text{large}} \approx 0$ .

→ So,



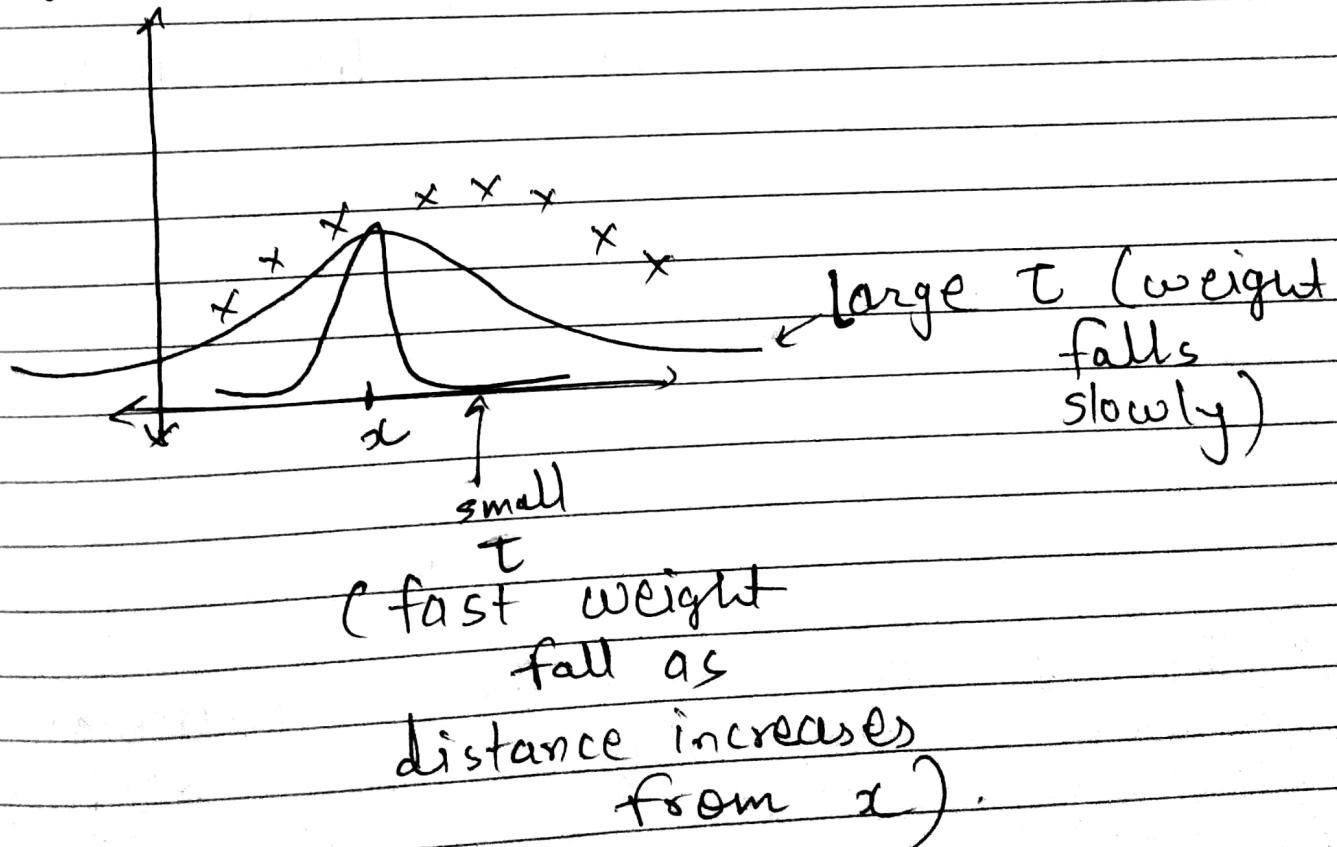
→ So, distant points from  $x$  does not contribute much to the above formula because  $\omega^{(i)}$  is near to zero.

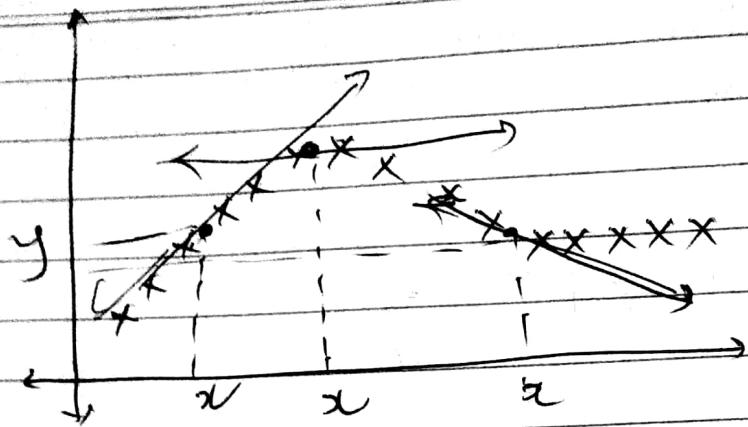
→ More general form of  $\omega^{(i)}$  can be given by

$$\omega^{(i)} = \exp\left(-\frac{(x_i - x)^2}{2T^2}\right)$$

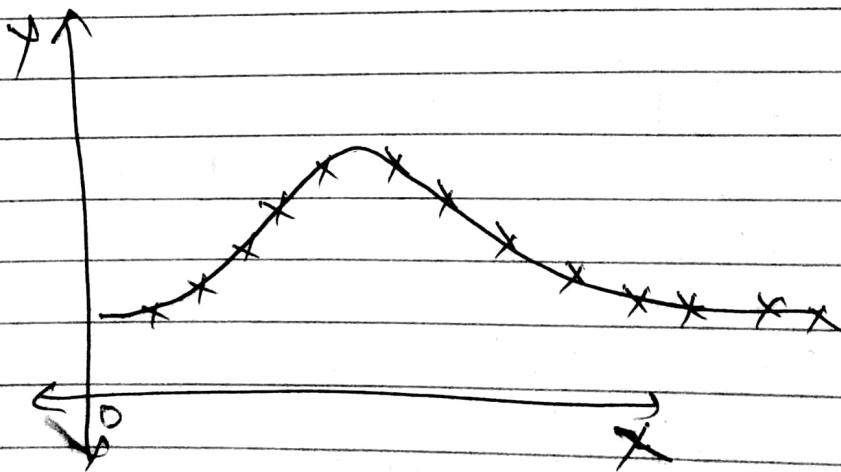
where  $T$  is bandwidth, which controls how fast weights of  $x_i$  fall with distance from  $x$ .

→ So, if  $T$  is small, then it will give narrow bell shape.





- so, for every new position prediction, we need to run fitting procedure at that point.
- So, if we do this at every point of  $x$ -axis then we get a full non-linear curve like following.



- Eq<sup>n</sup> of  $\omega^{(1)}$  looks like gaussian distribution, but it has nothing to do with it.

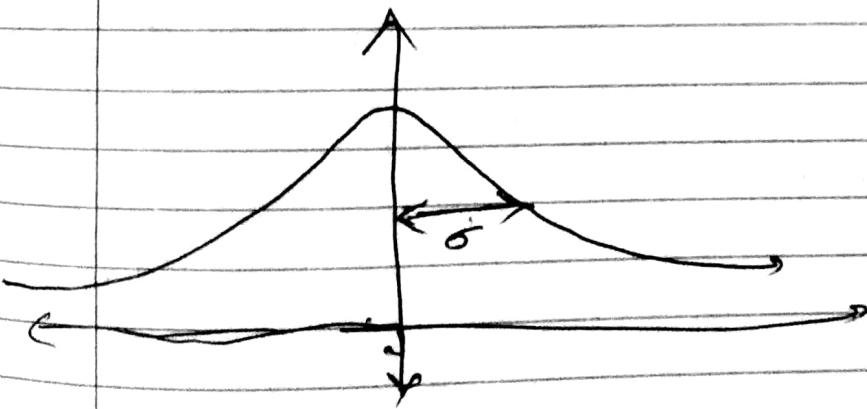
## \* Probabilistic Interpretation -

→ Let assume  $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$ ,

where  $\epsilon^{(i)}$  = error

$\epsilon^{(i)} \sim N(0, \sigma^2) \Rightarrow$  that means  $\epsilon^{(i)}$  are distributed as per gaussian distribution. And  $N$  is for Normal distribution.

$$P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right) \rightarrow ①$$



→ So, eqn ① means, that the probability distribution of the price of the house or ( $y^{(i)}$ ), for given  $x^{(i)}$  (features), and parameters  $\theta$ , will be gaussian

$$\text{So, } P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

or

$$y^{(i)} / x^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$$

→ here  $\theta$  is not random variable  
but  $\epsilon$  is.

→ Now, assume that error terms  $\epsilon^{(i)}$ 's  
are IID (independently and  
identically distributed).

→ So,

$$L(\theta) = P(\bar{y} | x; \theta) = \prod_{i=1}^m P(y^{(i)} / x^{(i)}; \theta)$$

↓  
 Likely hood  
 of  $\theta$ .  
 ↓  
 probability  
 of  $y$  given  $x$ ,  
 parameterized  
 by  $\theta$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

1

→ We need to choose  $\theta$  (parameters)  
to maximize likely hood  $L(\theta)$ .

\* Maximum likelihood :-

choose  $\theta$  to maximize  $L(\theta)$   
 $= P(\bar{y} | x; \theta)$

$$\begin{aligned}
 \rightarrow l(\theta) &= \log L(\theta) \quad \xrightarrow{\text{same as } (1)} \\
 &= \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp(-\dots) \\
 &= \sum_{i=1}^m \log \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp(-\dots) \right] \\
 &= m \log \frac{1}{\sqrt{2\pi}\sigma} + \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}
 \end{aligned}$$

$\rightarrow$  So, maximizing  $l(\theta)$  is the same as minimizing

$$\sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} = J(\theta)$$

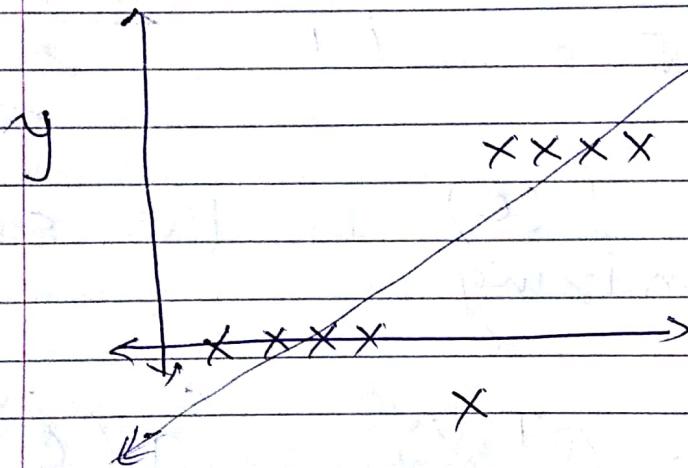
$\rightarrow$  Value of  $\sigma^2$  or other constants doesn't matter in minimizing the value of  $J(\theta)$ .

## (Bernoulli Distribution)

\* Classification Algo- In this,

the value of  $y$  or outcome is discrete values instead of continuous values as it was in regression problems.

→ for ex.  $y \in \{0, 1\}$ . or cancer type like (benign or malignant) etc.



If we fit linear line.  
but it won't be perfect.

→ So, if  $y \in \{0, 1\}$ , then we can define hypothesis as

$$h_0(x) \in [0, 1]$$

→ So, now instead of choosing linear function, we should choose non-linear function.

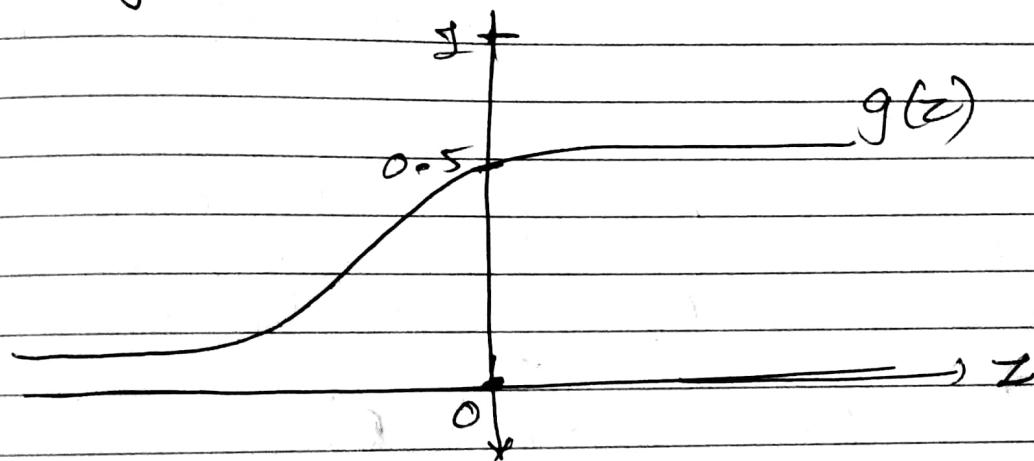
So,

$$h_0(x) = g(\theta^T x)$$

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

$$\therefore h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$$

$\rightarrow g$  is called sigmoid function or logistic function.



$\rightarrow$  Now,  $P(y=1 | x; \theta) = h_{\theta}(x)$ , then

$$P(y=0 | x; \theta) = 1 - h_{\theta}(x)$$

$\rightarrow$  Combining both of results

$$P(y|x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}$$

like binomial distribution

$$\begin{aligned} \rightarrow \text{So, } L(\theta) &= P(y|x; \theta) = \prod_i P(y^{(i)}|x^{(i)}; \theta) \\ &= \prod_i h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}} \end{aligned}$$

→ Now, to maximize the likelihood of  $\theta$ , let's assume

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

$$= \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)}))$$

→ Now, we can apply gradient descent to maximize the value of  $l(\theta)$  or likely hood. So,

$$\theta = \theta + \alpha \nabla_\theta l(\theta)$$

↓  
this is '+' here because  
we are maximizing  $\theta$ . It was  
'-' in gradient descent, bcz  
we were minimizing error.

→ So, this algo. can be called gradient Ascent.

$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

→ So,

$$\theta_j = \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

→ So, the formula is exactly same as batch gradient descent. But the only difference is value of  $h_\theta(x)$ .

→ here  $h_\theta(x)$  is  $\frac{1}{1+e^{-\theta^T x}}$  where as

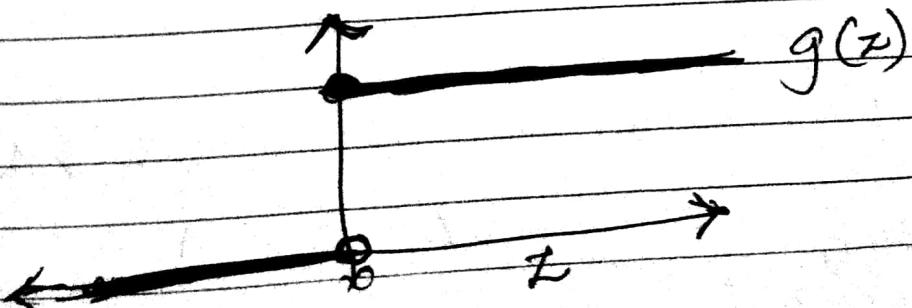
~~is~~ ~~in~~ ~~linear~~ ~~in~~ in gradient descent,

$$h_\theta(x) = \theta_0 + \theta_1 x_1$$

## \* Perceptron Algo

→ In sigmoid ~~func~~ function  $g(z)$  can give any value between  $[0, 1]$ . but if we only want to get value as 0 or 1, then it can be defined as

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



→ so, here  $h_{\theta}(x) = g(\theta^T x)$  (Same as linear regression)

&

$$\hat{Q}_j = Q_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x_j^{(i)}$$

also, same as previous logistic regression.

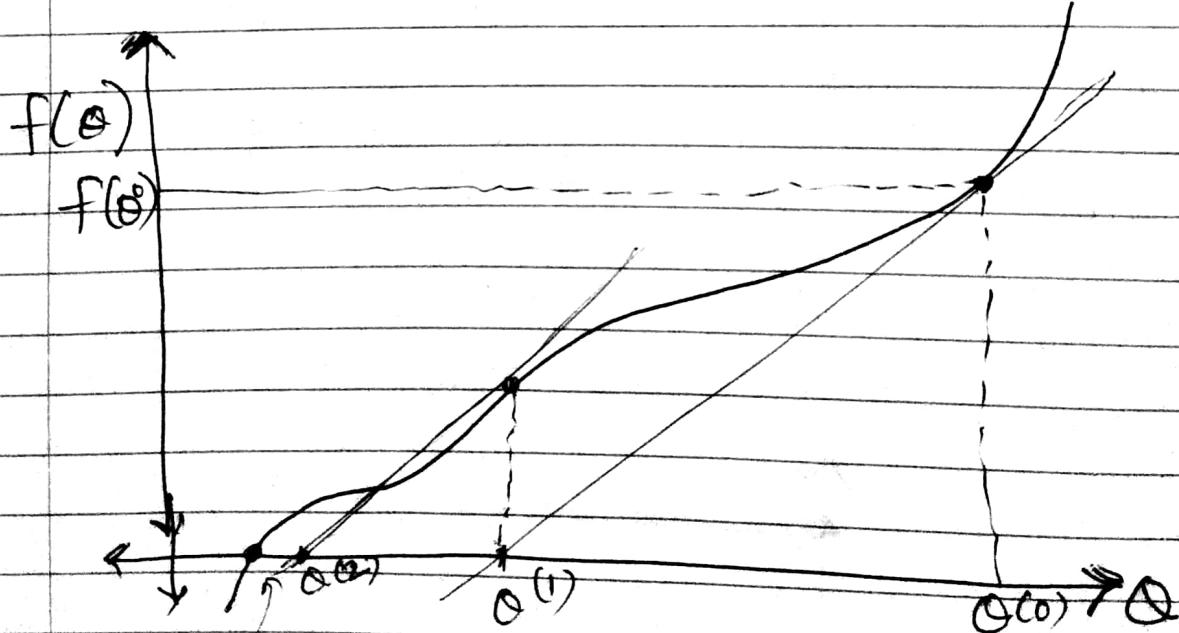
→ Perceptron algo is more simpler than logistic regression. (previous one).

Lecture-4 25/4/18

### \* Newton's Method:

→ Let's say we've a function  $f(\theta)$  and we need to find the value of  $\theta$ , such that  $f(\theta) = 0$ .

$f(\theta)$



we need to  
find this

initial guess

→ So, in this method, start with random guess of  $\theta$ , and then find slope at that point, now where the line (tangent at  $\theta^{(0)}$ ) intersects to  ~~$\theta$ -axis~~, the value at that point will be our next approx. of  $\theta$ .

$$\rightarrow f'(\theta^{(0)}) : \text{slop of tangent at } f(\theta^{(0)}) \\ = \frac{f(\theta^{(0)})}{\Delta}$$

$$\therefore \Delta = \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

$$\rightarrow \text{now, } \theta^{(1)} = \theta_0 - \Delta \quad (\text{from figure})$$

$$\theta^{(1)} = \theta_0 - \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

? In general,

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

→ one iteration of newton's method.

→ Now, this algo. can be used to maximize the log likelihood of  $\theta$  or  $[l(\theta)]$ .

→ So, we want  $\theta$ , such that  $l'(\theta) = 0$ , because  $l(\theta)$  is maximum, when  $l'(\theta) = 0$ , [Global Maxima].

so,

$$\theta^{(t+1)} = \theta^{(t)} - \frac{l'(\theta^{(t)})}{l''(\theta^{(t)})}$$

→ Here  $\theta$  is a ~~Real~~ Real number [ $\theta \in \mathbb{R}$ ], but if we've  $\theta$  as a vector [matrix] then the above generalized form will be

$$\theta^{(t+1)} = \theta^{(t)} - H^{-1} \nabla l,$$

where  $H^{-1}$  is a Hessian Matrix and  $\nabla l$  is gradient, as described previously.

$\Rightarrow$  In Hessian matrix,

$$H_{ij} = \frac{\partial^2 l}{\partial \theta_i \partial \theta_j}$$

A convergent rate is very high

→ Convergent Rate of Newton's method is very high as compared to gradient descent -

→ Convergent rate of Newton's method is **Quadratic**. So, if initial error is, ( $\text{error in } 1^{\text{st}} \text{ iteration}$ )

→ But the disadvantage of Newton's method is that we need to invert the Hessian Matrix in every iteration. Hessian matrix has dimension of  $(n+1) \times (n+1)$ , where  $n$  is no. of features.

## \* Generalized Linear Model :- (GLM)

→ Suppose we've the data that is discrete ( $\{0, 1\}$  valued) and we want to model it with Bernoulli Random variable, parameterised by  $\phi$ ,

Bernoulli( $\phi$ ),

→ it's not a single distribution, but its set of distributions

$$\rightarrow P(y=1; \phi) = \phi$$

→ As we vary the value of  $\phi$ , we'll get different probability distributions.

→ Bernoulli and Gaussian distributions are special cases of exponential family of distribution.

$$\rightarrow P(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

where,

$\eta$  = natural param of distribution  
 $T(y)$  = sufficient statistic

(Usually  $T(y) = y$ ).

- Now let's show that bernoulli and Gaussian distribution is special case of exponential family of distribution.
- Let's start with Bernoulli.

→  $Ber(\phi)$ , ~~if~~  $P(y=1; \phi) = \phi$

- Now, let's choose  $a, T$  and  $b$ , such that formula (1) becomes identical to the above one.

$$\begin{aligned}
 \rightarrow \text{so, } P(y; \phi) &= \phi^y (1-\phi)^{1-y} \\
 &= \exp(\log \phi^y (1-\phi)^{1-y}) \\
 &= \exp(y \log \phi + (1-y) \log(1-\phi)) \\
 &= \exp(y \log \phi + \log(1-\phi) - y \log(1-\phi)) \\
 &= \exp(y \underbrace{\log \phi}_{T(y)} + \underbrace{\log(1-\phi)}_{-a(\eta)})
 \end{aligned}$$

$$P(y; \phi) = \exp \left( \underbrace{y \log \phi}_{T(y)} + \underbrace{\log(1-\phi)}_{-a(\eta)} \right)$$

$$\text{and } b(y) = 1.$$

$$\therefore \eta = \log \frac{\phi}{1-\phi}$$

$\therefore \phi = \frac{1}{1+e^{-\eta}}$   $\Rightarrow$  Same as logistic function or sigmoid f<sup>n</sup>.

$$\rightarrow a(\eta) = -\log(1-\phi)$$

$$a(\eta) = \log(1+e^\eta)$$

$$\rightarrow T(y) = y \quad \text{and} \quad b(y) = 1$$

$\rightarrow$  Now let's do the same for Gaussian.

~~33:22~~  $N(\mu, \sigma^2)$  set  $\sigma^2 = 1$ , because  $\sigma^2$  value doesn't matter in maximizing of likelihood of  $\ell(\theta)$ .

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

$\rightarrow$  taking this in terms of formula (1)  
we will have

$$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \exp\left(\frac{w-y}{\sqrt{2}} + \frac{1}{2}w^2\right)$$

$b(y)$        $w=\mu$        $T(y)=y$        $= \frac{1}{2}w^2$

\* Having chosen exponential family distribution, how to derive a generalized linear model? ~~GLM~~ (GLM) \*

→ Assume:  $y|x; \theta$  is distributed.

1)  $y|x; \theta \sim \text{Exp Family}(y)$

2) Given  $x$ , our goal is to output expected value of  $y$ , given  $x$ .  
i.e.  $E[y|x] = E[T(y)|x]$

In other words,  $h(x) = E[T(y)|x]$   
(usually  $T(y) = y$ ).

3)  $\eta = \theta^T x$  [design choice]  
this is when  $y$  is real number.  
but if  $y$  is a vector, then

$$[y_i = \theta_i^T x, \text{ if } \eta \in \mathbb{R}^k]$$

→ For Bernoulli distribution,

$y|x; \theta \sim \text{Exp Family}(y)$

→ for any fixed value of  $x$  and  $\theta$ , our learning algo, make prediction,

$$h_\theta(x) = E[y|x; \theta] = P(y=1|x; \theta) = \phi$$

$$= \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}}$$

$$\rightarrow g(\eta) = E[y; \eta] = \frac{1}{1+e^{-\eta}}$$

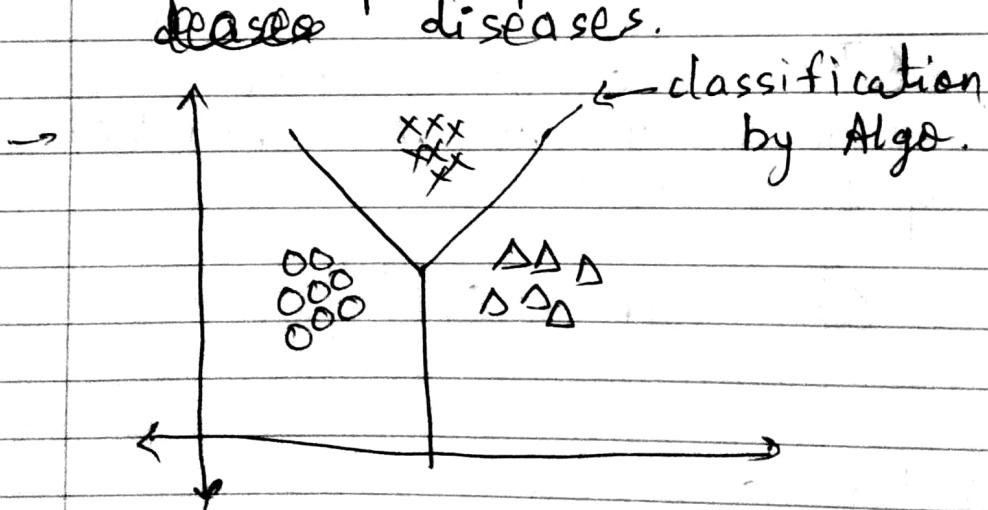
which is canonical response  $f^n$ .  
 $g^{-1}(\eta)$  is canonical link function.

### \* Multinomial :-

→ It's a distribution over  $k$  discrete values  
 So,

$$y \in \{1, 2, \dots, k\}$$

→ for example, patient has which of  $k$   
 classes diseases.



→ So, parameters:-  $\phi_1, \phi_2, \phi_3, \dots, \phi_k$

where,  $P(y=i) = \phi_i$

$$\phi_1 + \phi_2 + \dots + \phi_k = 1$$

$$\therefore \phi_k = 1 - (\phi_1 + \phi_2 + \dots + \phi_{k-1})$$

- for this derivation, let take  $\phi_1, \phi_2, \dots, \phi_{k-1}$  as parameters.  
and  $\phi_k = 1 - (\phi_1 + \phi_2 + \dots + \phi_{k-1})$ .
- here  $y \in \{1, 2, \dots, k\}$ .

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \in \mathbb{R}^{k-1}$$

$$T(k-1) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \quad \text{and} \quad T(k) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{k-1}$$

→ Note that here  $T(y) \neq y$

→ Consider,

$$1\{\text{TRUE}\} = 1, \quad 1\{\text{FALSE}\} = 0$$

so, if  $1\{z=3\} = 0$  and  $1\{1+1=2\} = 1$

→  $T(y)_i = 1\{y=i\} \Rightarrow$  this means,  
that the  $i^{\text{th}}$  element of vector  $T(y)$   
will only be equal to 1, iff  
 $y = i$ .

$$\rightarrow P(y) = \phi_1^{1(y=1)} \phi_2^{1(y=2)} \cdots \phi_k^{1(y=k)}$$

Remember  
this is not  
parameter

$$= \phi_1^{T(y)_1} \phi_2^{T(y)_2} \cdots \phi_{k-1}^{T(y)_{k-1}} \phi_k^{1 - \sum_{j=1}^{k-1} T(y)_j}$$

= ... [after few steps]

~~$\rightarrow = b(y) \exp(\eta^T T(y) - a(\eta))$~~

where,

$$\eta = \begin{bmatrix} \log(\phi_1/\phi_k) \\ \vdots \\ \log(\phi_{k-1}/\phi_k) \end{bmatrix} \in \mathbb{R}^{k-1}$$

*(making subject)*  
 $\rightarrow a(y) = -\log(\phi_k), \text{ and } b(y) = 1.$

$$\phi_i = \frac{e^{n_i}}{1 + \sum_{j=1}^{k-1} e^{n_j}}, (i=1, \dots, k-1)$$

→ using the assumption,

$$\phi_i = \frac{e^{\phi_i^T x}}{1 + \sum_{j=1}^{k-1} e^{\phi_j^T x}}, [n_i = \phi_i^T x]$$

$$\rightarrow h_{\theta}(x) = E[T(y) | x; \theta]$$

$$= E \left[ \begin{array}{c} 1\{y=1\} \\ \vdots \\ 1\{y=k-1\} \end{array} \mid x; \theta \right]$$

$$= \begin{bmatrix} \phi_1 \\ \phi_2 \\ \phi_3 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

$\rightarrow$  This algo. is called Softmax Regression  
 [Generalization of Logistic Regression].

$\rightarrow$  In a particular machine learning problem, if,

$y \in \{1, \dots, k\}$  [ex.  
 Prediction of digits from images].

$\Leftrightarrow$  our training set will be

$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ , where

$y^{(i)}$  can take any value from  $\{1, \dots, k\}$ ,  
 then, we'll try to maximize the likelihood of  $\theta$ , as following.,

$$L(\theta) = \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta).$$

$$= \prod_{i=1}^m \phi_1^{1\{y^{(i)}=1\}} \cdot \phi_2^{1\{y^{(i)}=2\}} \cdots \phi_k^{1\{y^{(i)}=k\}}$$

where  $\phi_i = \frac{e^{\theta_i^T x}}{1 + \sum_{j=1}^k e^{\theta_j^T x}}$

→ here  $\theta_1, \theta_2, \dots, \theta_k \in \mathbb{R}^{n+1}$ , where n is no. of features.

## Lecture-5

### \* Generative learning Algorithms

- In this we train different model for each class and then compare the data with that models to predict the result.
- For example, to predict cancer, whether it is malignant or benign we train separate model for both & then at the time of prediction we compare our data to both the model.

- \* Logistic regression is discriminative learning algo. so,
- it learns  $P(y|x)$  or learns  $h_0(x) \in \{0, 1\}$  directly.
- Generative algo. builds the model on what the features looks like conditioned on class label. so,

$$P(x|y) \text{ and } P(y)$$

$$P(y=1|x) = \frac{P(x|y=1) P(y)}{P(x)}$$

$$\therefore P(x) = P(y=0|x) P(x) + P(y=1|x) P(x)$$

Ex Assume that the input features  $x \in \mathbb{R}^n$  and continuous valued.

→ Gaussian Discriminant Analysis  
Here assume that  $P(x|y)$  is Gaussian.

→ Multi variable Gaussian ~~can~~ can be given by,  $\xrightarrow{\text{mean matrix}} \Sigma \xrightarrow{\text{covariance matrix}}$

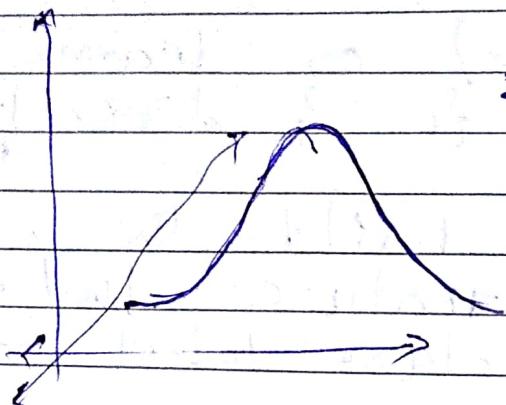
$$x \sim N(\bar{u}, \Sigma)$$

$$P(x) = \frac{1}{\sqrt{2\pi}^n |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2} (x-\bar{u})^T \Sigma^{-1} (x-\bar{u}))$$

$$\text{where } \Sigma = E[(x-\bar{u})(x-\bar{u})^T]$$

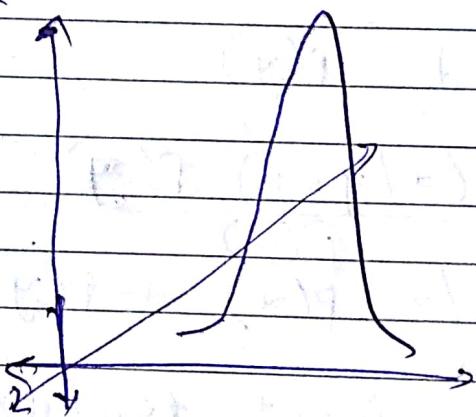
→ Examples of Gaussian with mean 0 and different co-variance matrix.

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



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

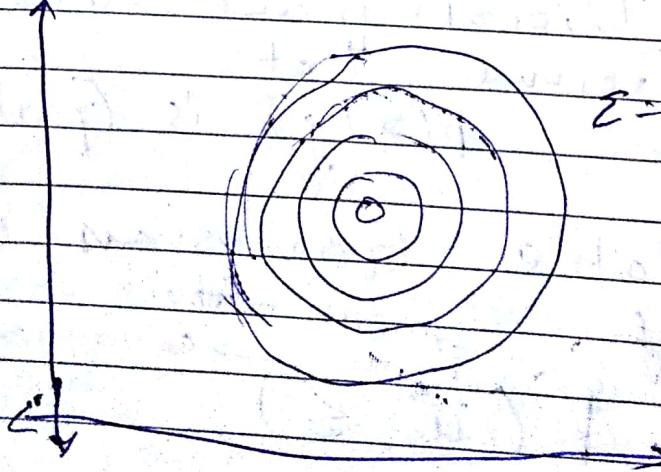
↑ responsible  
for shape of curve.  
↓ responsible  
for location

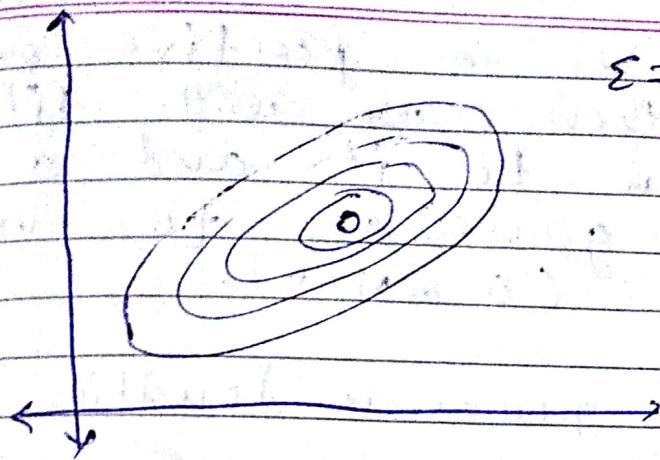


$$\Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$

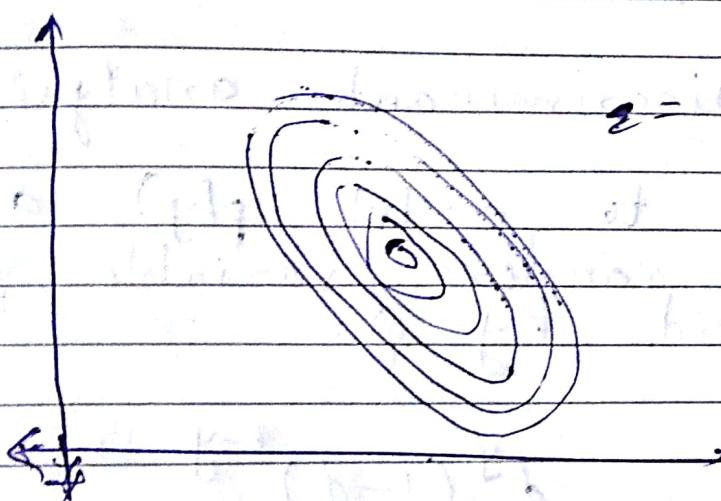
→ In 2-D it can be shown as. (Top view)

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



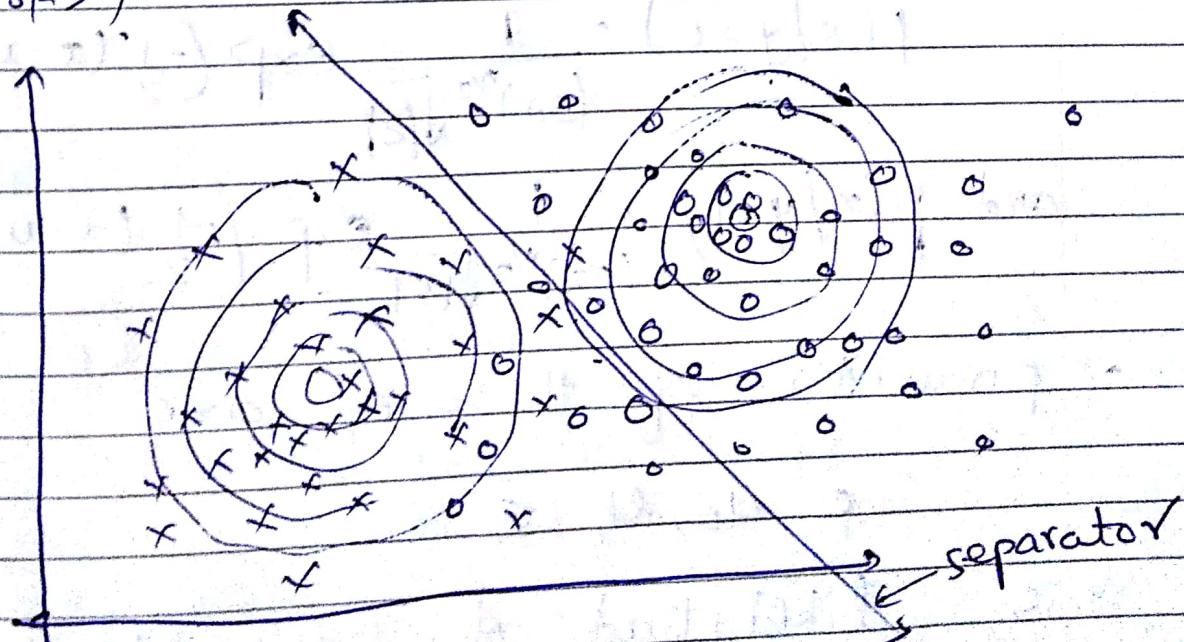


$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

\* How Gaussian discriminant analysis works?



we will fix a gaussian to  $x$  and different gaussian to  $o$

- & let's say 'x' is positive examples ( $y=1$ ) then we will fit a gaussian to it and a different gaussian to -ve examples ('0' or  $y=0$ ).
- And this two gaussian densities defines a separator.
- & Gaussian discriminant analysis model-
- It's going to model  $p(y)$  as a bernoulli random variable, ~~parameterized~~ parameterized by  $\phi$ .

so,

$$p(y) = \phi^y (1-\phi)^{1-y}$$

$$p(x|y=0) = \frac{1}{(2\pi)^{n/2} \sqrt{|S|}} \exp\left(-\frac{1}{2} (x-u_0)^T S^{-1} (x-u_0)\right)$$

$$\text{and } p(x|y=1) = \frac{1}{(2\pi)^{n/2} \sqrt{|S|}} \exp\left(-\frac{1}{2} (x-u_1)^T S^{-1} (x-u_1)\right)$$

→ parameters of these eq<sup>n</sup> are

$$\phi, u_0, u_1, S$$

→ So, likelihood of parameters

also called  
Joint likely hood.

Page No.	
Date	

$$l(\phi, u_0, u_1, \Sigma) = \log \prod_{i=1}^m P(x^{(i)}, y^{(i)})$$

$$= \log \prod_{i=1}^m P(x^{(i)}|y^{(i)}) \cdot P(y^{(i)})$$

→ Now, as previous, to maximize  $l$ ,  
w.r.t  $\phi$ ,  $u_0$ ,  $u_1$  and  $\Sigma$ ,

where,

$$\text{max. likely } \phi \underset{\text{hood}}{=} \frac{\sum_i y^{(i)}}{m} = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=1\}}{m}$$

$$u_0 = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=0\} x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=0\}}$$

$$u_1 = \frac{\sum_{i=1}^m (\mathbb{1}\{y^{(i)}=1\} x^{(i)})}{\sum_{i=1}^m \mathbb{1}\{y^{(i)}=1\}}$$

no. of examples with  
label 1.

indicates  
only the terms  
where  $y^{(i)}=1$ .

so, numerator says sum of  $x^{(i)}$

corresponding to the examples,  
where  $y^{(i)} = 1$ .

→ it simply indicates that-  
find all the examples for  
 $y=1$  (+ve examples) and  
take the ~~avg.~~ average

→ predict:

$$\arg \max_y P(y|x) = \arg \max_y \frac{P(x|y)P(y)}{P(x)}$$

$$= \arg \max_y P(x|y) P(y)$$

$P(x)$  is ignored because it  
does not depend on  $y$ .

→ If  $P(y)$  is uniform, ( $P(y=0) = P(y=1)$ )

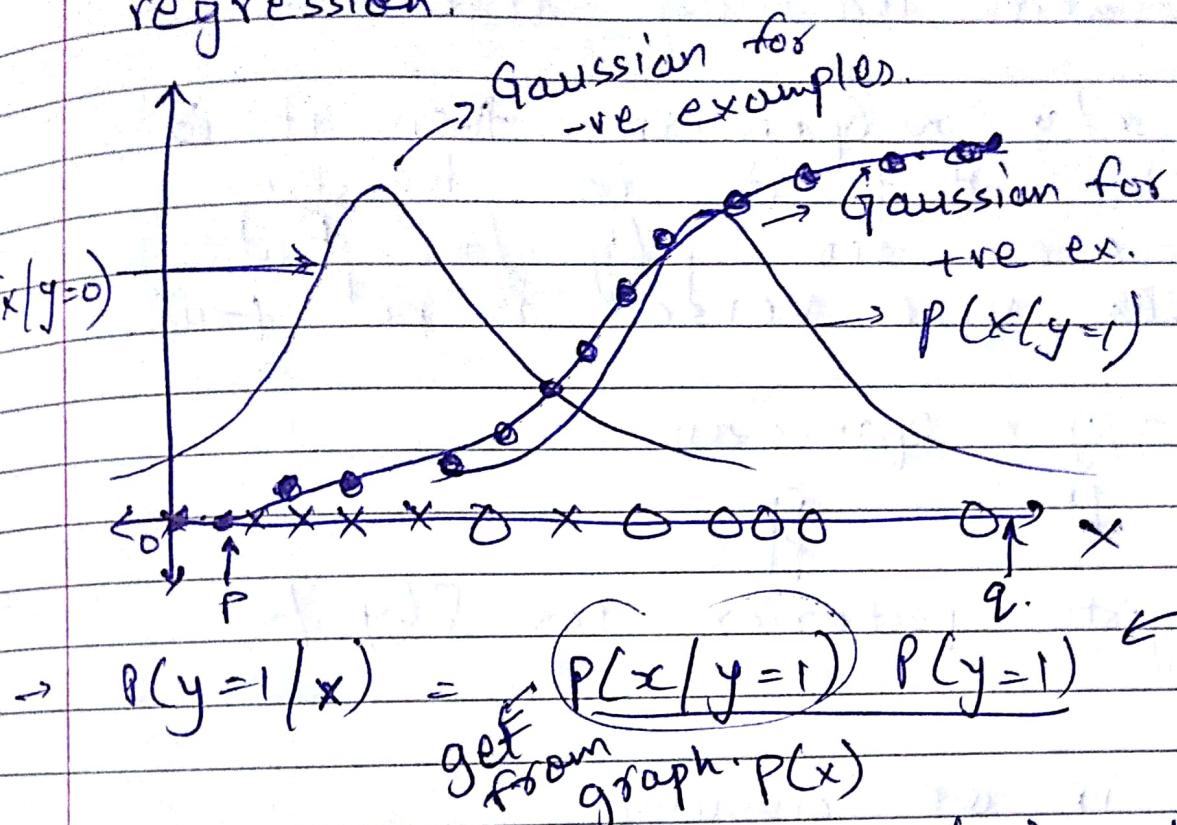
then,

$$= \arg \max_y P(x|y)$$

\* Note-  $\arg \max_y P(x|y)$  means that

the value of  $y$ , @ for which  
 $P(x|y)$  is maximum.

→ Gaussian discriminant analysis has interesting relation with logistic regression.



$$\rightarrow P(y=1|x) = \frac{P(x|y=1) P(y=1)}{\text{get from graph } P(x)} \quad \text{this is simply a Bayes rule.}$$

→ If we calculate  $P(y=1|x)$  at point p, then it will be very low. (almost 0)

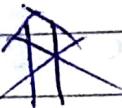
→ As we keep moving right side on x-axis, the probability increases. If we calculate  $P(y=1|x)$  at point q, it will be almost 1. So, If we plot the curve for all such points, it will give the curve of logistic regression.

$$\Rightarrow \text{and } P(x) = P(x|y=0) P(y=0) + P(x|y=1) P(y=1)$$

## \* Advantages and Disadvantages of Generative learning algo.

- if  $x/y \sim \text{Gaussian}$ , then it ~~is~~ implies that it is logistic posterior for  $p(y=1/x)$ . But vice versa is not true

$x/y \sim \text{Gaussian}$

$\therefore$  

Logistic posterior for  $P(y=1/x)$

- Also, if we assume,

$$\begin{cases} x/y=1 \sim \text{Poisson}(\lambda_1) \\ x/y=0 \sim \text{Poisson}(\lambda_0) \end{cases}$$

then  $P(y=1/x)$  is logistic.

- So,  $x/y$  is Gaussian is ~~a~~ stronger assumption ~~as~~, as, it implies if  $y/x$  is logistic but vice versa is not true.

→ Generative learning Algo. used when data set is small. And also gives good results ~~if~~ when such assumptions do not hold.

→ More generally,

$$\begin{aligned} x/y = 1 &\sim \text{Exp Family } (n_y) \\ x/y = 0 &\sim \text{Exp Family } (n_0) \end{aligned}$$

If implies

$P(y=1/x)$  is logistic.

→ Other Generative learning algo. is Naive Bayes Algo.

→ For ex. we want to classify spam emails, then we need to represent an email mathematically. So, if we create a matrix (vector) of dictionary, then email can be represented by as following.

[	1	a
0	aardvark	
0	:	
:		
0	buy	
0	:	
:		

problem is in  
the form of  
 $y \in \{0, 1\}$   
not  
spam      spam

→ now, we want to find  $P(x|y)$ ,  
where  $x \in \{0, 1\}^n$

$n = 50,000$  or more dimensions.

→ so, if we use softmax or multinomial regression here,  
then no. of possible matrix are  $2^{50,000}$ , so,  $2^{50,000}-1$  params.

→ Assume  $x_i$ 's are conditionally independent given  $y$ .

$$P(x_1, x_2, x_3, \dots, x_{50,000} | y) = P(x_1 | y) \cdot P(x_2 | y, x_1) \cdot P(x_3 | y, x_1, x_2) \cdots$$

→ also assume in above eq<sup>n</sup>, that

$$P(x_2 | y, x_1) = P(x_2 | y)$$

$$P(x_3 | y, x_1, x_2) = P(x_3 | y) \cdots$$

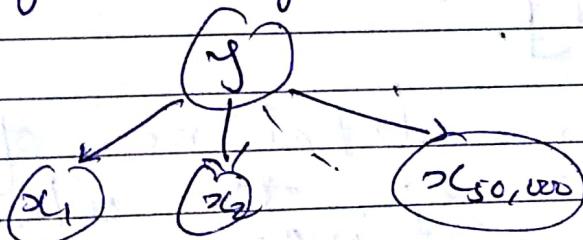
$$\text{so, } P(x_1, x_2, \dots | y) = \prod_{i=1}^n P(x_i | y)$$

\* What  $P(x_2 | y, x_1)$  means?

→ It means, given the value of  $y$  (0/1) [spam or not spam], and given the word  $x_1$  appears in email, does not affect the probability

of word  $x_2$  appears in email.

- But practically this assumption is false, for example if there is word 'CS229' in email, then word 'Andrew Ng' will be there. But despite of this false assumption, the algo. works correctly in practice.
- Particularly, this algo. works very well for classifying text archives.
- Bayesian Network for this model can be given by,



- Following are the parameters of this models.

$$\phi_{i|y=1} = P(x_i=1 | y=1) \quad | \quad \text{if } y=1 \text{ means}$$

$$\phi_{i|y=0} = P(x_i=1 | y=0) \quad | \quad \text{that the fraction of spam emails}$$

$$\phi_y = P(y=1)$$

$(y=1)$  in which word  $i$  does appear.

- Joint likelihood:

$$f(\phi_y, \phi_{j|y=1}, \phi_{i|y=1}) = \prod_{i=1}^m TTP(x^{(i)}, y^{(i)})$$

→ we can get value of this params based on max. likelihood analysis.

$$\phi_{j|y=1} = \frac{\sum_{i=1}^m I\{x_j^{(i)}=1, y^{(i)}=1\}}{\sum_{i=1}^m I\{y^{(i)}=1\}}$$

→ This means:-

denominator ← Total no. of spam emails in training set. [emails where  $y=1$ ].

numerator ← total no. of words  $x_j$ , for which the email was spam.

$$\phi_y = \frac{\sum_{i=1}^m I\{y^{(i)}=1\}}{m}$$

→ To predict for a new email, compute  $P(y|x)$  using bayes rule.

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

## \* Laplace smoothing:

- For Example, if any  $w$  is in email classification problem using naive bayes algo, if we see a new word 'nips' which was never seen before in any of the training set.
- Assume that this word is 35000<sup>th</sup> word in the dictionary.
- So, our spam classifier algo. will predict following probabilities.

$$\phi_{35000/y=1} = \frac{\sum_{i=1}^m I\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m I\{y^{(i)} = 1\}} = 0$$

and

$$\phi_{35000/y=0} = \frac{\sum_{i=1}^m I\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m I\{y^{(i)} = 0\}} = 0$$

- This is bcz it has never seen 'nips' in either spam or non spam. So it thinks seeing this in either type of email has probability 0.

→ And from this, when it computes  
(using bayes rule)  
~~product~~

$$\begin{aligned} P(y=1/x) &= \frac{P(x/y=1) P(y=1)}{P(x)} \\ &= \frac{P(x/y=1) P(y=1)}{P(x/y=1) P(y=1) + P(x/y=0) P(y=0)} \\ &\stackrel{50,000}{\overbrace{\prod_{i=1}^{50,000} P(x_i/y=1)}} \end{aligned}$$

one of the term in this product will be zero because of previous calculation.

$$= \frac{0}{0+0} = \frac{0}{0} \text{ (undefined).}$$

→ But the problem is that actually  $P(x_{35000}/y=1) = 0$  is bad assumption.

→ We can't set probability of this to zero, as we haven't seen the word 'nips' before.

→ So we need to fix this problem.

\* Example of basket ball team :

Date	vs match	result (0=lose, 1=win)
2/8	Washington state	0
2/11	Washington	0
2/22	USC	0
2/24	UCLA	0
3/8	USC	0
3/15	Louisville	?

→ So, here we can't directly predict that there is no chance to win the last match just because of previous data.

→ Here comes the Laplace smoothing.

→ Normally the max likelihood of ( $y=1$ ) can be given by (informally) -

$$P(y=1) = \frac{(\# \text{ of } 1)}{(\# \text{ of } 0) + (\# \text{ of } 1)}$$

total training

data

→ Laplace just add 1 to each term above

$$\text{So, } P(y=1) = \frac{(\# \text{ of } 1) + 1}{(\# \text{ of } 1) + 1 + (\# \text{ of } 0) + 1}$$

→ So, in our ex.

$$\begin{aligned} P(y=1) &= \frac{1}{5+2+1} \\ &= \frac{1}{7} \quad \text{instead of } 0. \end{aligned}$$

probability  
of winning  
next game

→ More generally,

if  $y \in \{1, 2, \dots, k\}$  (instead of 0 & 1 as above)

$$P(y=j) = \frac{\sum_{i=1}^m I\{y^{(i)} = j\}}{m+k} + 1$$

→ so for naive bayes,

$$\begin{aligned} \phi_j | y=1 &= \frac{\sum_{i=1}^m I\{x_j^{(i)} = 1, y^{(i)} = 1\}}{\sum_{i=1}^m I\{y^{(i)} = 1\}} + 1 \\ &= \frac{1}{m+2}. \end{aligned}$$

and,

$$\phi_{j|y=0} = \sum_{i=1}^m I\{x_j^{(i)} = 1, y^{(i)} = 0\} + 1$$

$$\sum_{i=1}^m I\{y^{(i)} = 0\} + 2$$

## Lecture-6

\* Multinomial Event model for Naive Bayes.

→ In this model, representation of email will be different.

(i) is  $\rightarrow (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})$ , where  $n_i = \# \text{ of words}$   
 in an email  
 training example.

→ So, there email is represented by feature vector.

$x_j^{(i)} \in \{1, 2, \dots, \text{size}\}$ , so  $j$  is index of word  $x_j^{(i)}$  in dictionary

→ So, joint distribution is given by

$$P(x, y) = \left( \prod_{i=1}^n P(x_i|y) \right) P(y)$$

where  $n$  is length of an email

→ meaning of this  $\phi_k^n$  is that as it is generative learning algo, first  $P(y)$  is chosen [Email ~~is~~ is spam or not spam], ~~too~~ and then for that email, probability of all the words, depending on email is spam or not is calculated,

→ so, params of this,

$$\phi_{k|y=1} = P(x_j=k | y=1)$$

$$\phi_{k|y=0} = P(x_j=k | y=0)$$

$$\phi_y = P(y=1)$$

→ Max. likelyhood of param. can be given by

$$\phi_{k|y=1} = \sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\} \sum_{j=1}^{n_i} \mathbb{I}\{x_j^{(i)} = k\} + 1$$

$$\sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\} = \sum_{i=1}^m n_i$$

→ Numerator means, take all the spam emails, and in each such email, count the word  $k$ .   
 Laplace smoothing

- Denominator means that sum over entire training set, when email is spam, add the length of that email to the sum.
  - So, it is sum of length of all spam emails.
  - Note: If  $x$  can get ' $l$ ' different values, i.e.  $x \in \{1, 2, \dots, l\}$ , then, max. likelihood of  $x=k$  ( $k \leq l$ ) is given by
- $P(x=k) = \frac{\text{\# of observations of } x=k}{\text{total observations.}}$
- and with Laplace smoothing
- $P(x=k) = \frac{\text{\# of obs. when } x=k + 1}{\text{total obs.} + l}$
- How Maximum likelihood formula are derived?

$$\begin{aligned}
 L(\phi_{k|y=1}, \phi_{k|y=0}, \phi_y) \\
 &= \log \prod_{i=1}^m P(x^{(i)}, y^{(i)}; \phi_{k|y=1}, \phi_{k|y=0}, \phi_y)
 \end{aligned}$$

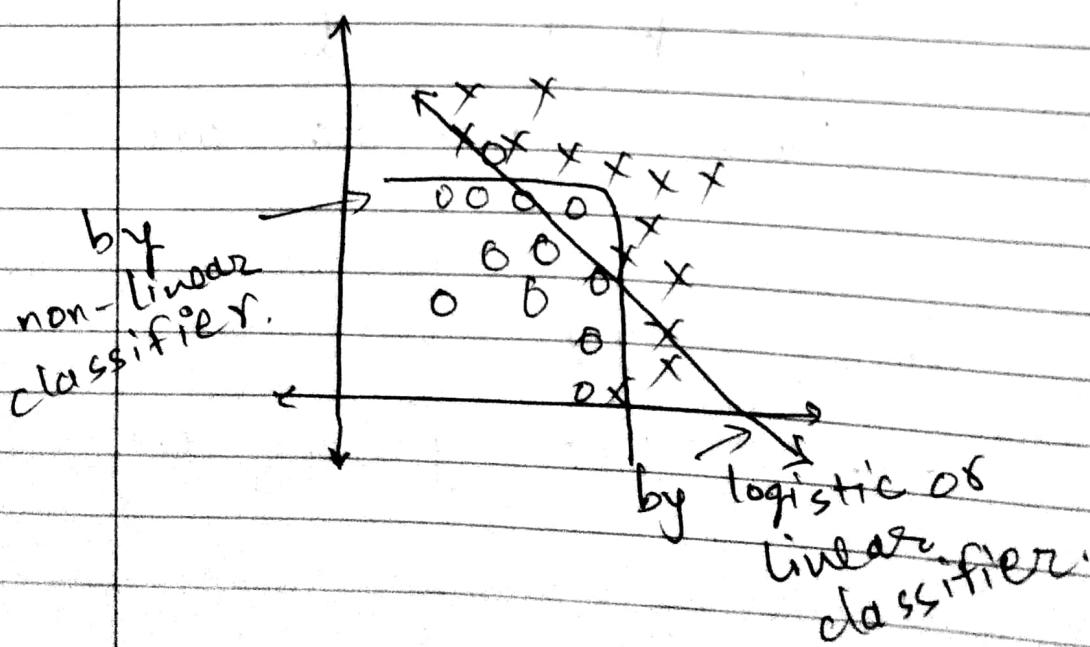
$$= \log \prod_{i=1}^m \prod_{j=1}^{n_i} P(x_j^{(i)} | y^{(i)}; \theta_k | y=1, \theta_k | y=0) \cdot P(y^{(i)}; \theta_y)$$

→ And by maximizing the above eq,  
 we can get likelihood estimation  
 of the parameters of the model.  
 [Maximizing means  $\frac{\partial}{\partial (\text{param})} = 0$   
 or any other algo like newton  
 raphson, gradient descent etc]

### \* Non-linear classifiers

→ The first classifier we've learned  
 was logistic regression which is  
 linear classifier

→ so, for ex. data is



→ Generative learning algorithms like Naive Bayes ~~or~~ and Gaussian discriminant Analysis ~~are~~ are also linear classifiers. Because under those models, we use linear classifiers as discussed in previous section of relation between GDA and Logistic Regression model.

So, as discussed before, if

$p(y=1|x) \sim \text{ExpFamily}(\eta_1)$   
and  $p(y=0|x) \sim \text{ExpFamily}(\eta_0)$

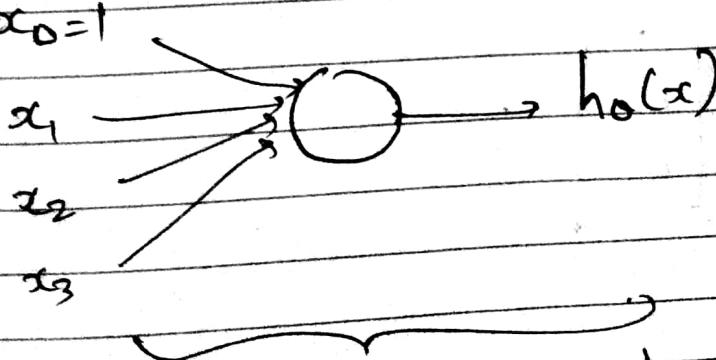
It implies that

$p(y=1|x)$  is logistic

$$\Leftrightarrow h_0(x) = \frac{1}{1 + e^{-\eta_0^T x}}$$

→ Consider a notation for further non-linear classification algorithms.

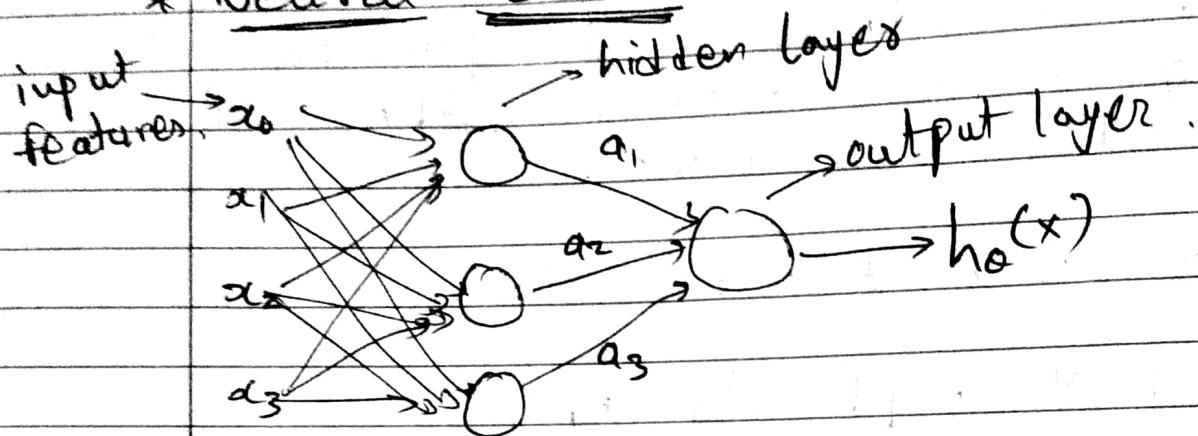
feature



Logistic Regression  
Unit:

→ One way to come up with a non-linear boundary, is Neural Network.

### \* Neural Network:



$a_1 = g(x^T \theta^{(1)})$  where  $g$  is sigmoid function and  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$  are different parameters associated with each logistic unit above.

$$h_0(x) = g(\vec{a}^T \theta^{(4)}), \text{ where}$$

$$\vec{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}.$$

$h_0(x)$  is function of  $\theta_1, \theta_2, \theta_3$  &  $\theta_4$ .

→ now, as before cost function is,

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2$$

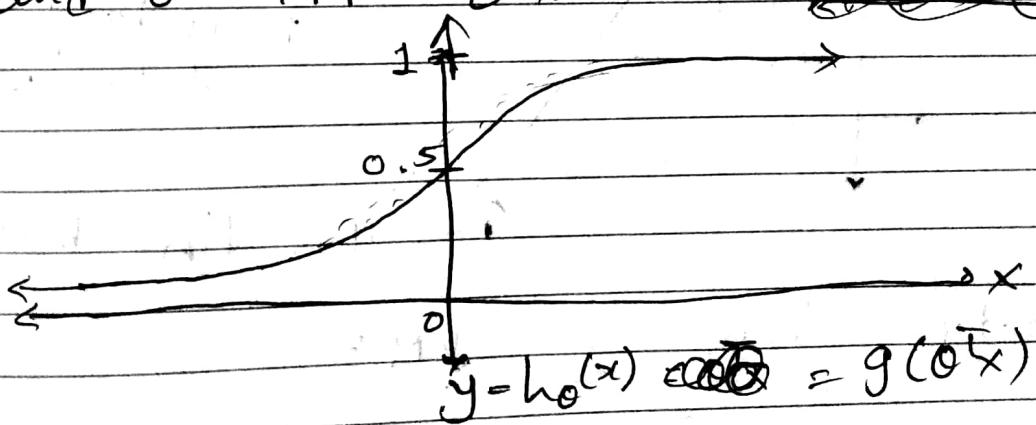
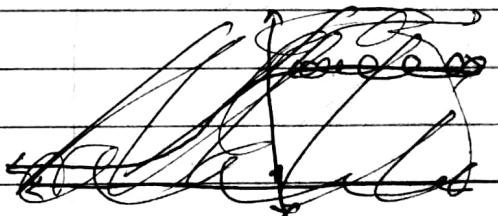
→ To learn the parameters ( $\theta_0, \theta_1, \dots, \theta_n$ ) we can use gradient descent.

→ Now Gradient Descent on Neural Network is called Back Propagation.

\* Support Vector Machines = (SVM). A supervised learning algo, classification algo.

→ In Logistic regression, we compute  $\theta^T x = y$  and predict

'1' iff  $\theta^T x > 0$   
and '0' iff  $\theta^T x < 0$



→ if  $\theta^T x \gg 0$ , very 'confident' prediction that  $y=1$ .  
if  $\theta^T x \ll 0$ , very 'confident' prediction that  $y=0$

→ It would be nice, if

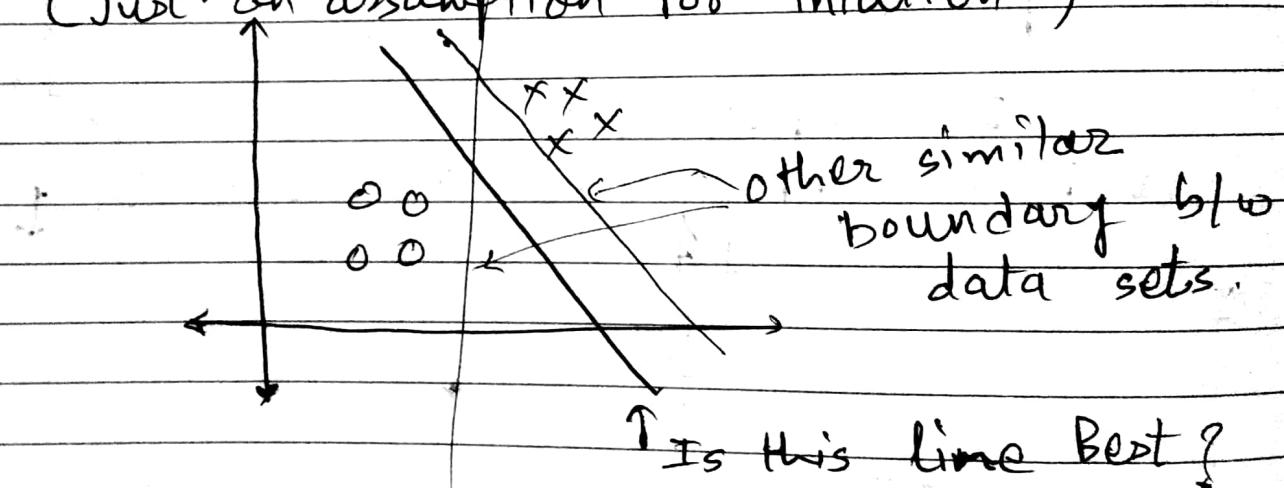
for  $\forall i$ , such that  $y^{(i)} = 1$ , we've  
 $\Omega^T x^{(i)} > 0$

and  $\forall i$ , such that  $y^{(i)} = 0$ , we've  
 $\Omega^T x^{(i)} < 0$ ,

→ So, our predictions be very confident.

→ The above intuition is called functional margin.

→ The 2<sup>nd</sup> intuition is called Geometric Margin. If dataset is linearly separable (Just an assumption for intuition).



→ Notation change of sum :-

→  $y \in \{-1, +1\}$  instead of  $y \in \{0, 1\}$

→ Have  $h$  (hypothesis) output values in  $\{-1, +1\}$

→  $g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{otherwise} \end{cases}$

$$\rightarrow h_0(x) = g(\omega^T x) \quad x \in \mathbb{R}^{n+1} \quad \left. \begin{array}{l} \uparrow \\ \omega_0 = 1 \end{array} \right\} \text{DROP IT}$$

Instead, of this,

$$h_{\omega, b}(x) = g(\omega^T x + b)$$

$\left[ \begin{array}{c} \omega_1 \\ \omega_2 \\ \vdots \end{array} \right]$

$\omega_0$

where  $\omega \in \mathbb{R}^n$  and  $x \in \mathbb{R}^n$  and  $b \in \mathbb{R}$

$\rightarrow$  Def<sup>n</sup>: Functional margin of a hyperplane  $(\omega, b)$  w.r.t training examples  $(x^{(i)}, y^{(i)})$  is defined as,

$$\hat{\gamma}^{(i)} = y^{(i)} (\omega^T x + b).$$

$\rightarrow$  Informally hyperplane means decision boundary.

If  $y^{(i)} = 1$ , we want  $\omega^T x + b \gg 0$   
 if  $y^{(i)} = -1$ , we want  $\omega^T x + b \ll 0$

so if  $y^{(i)} (\omega^T x^{(i)} + b) > 0$ , then we've classified  $(x^{(i)}, y^{(i)})$  correctly.

→ functional margin w.r.t the entire training set can be given by

$$\hat{\gamma} = \min_i |\hat{y}^{(i)}|$$

→ But for our case let's define that we want worst case functional margin. [Large fm margin]

→ To make large functional margin, we can just scale the parameters

$$w \rightarrow 2w$$

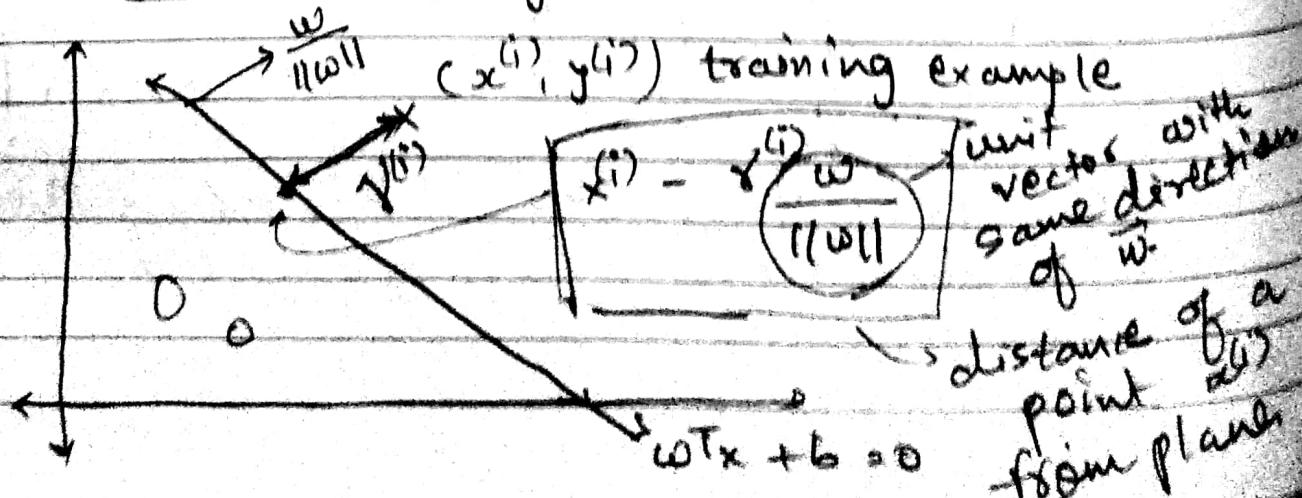
$$b \rightarrow 2b, \text{ then as per}$$

$\hat{y}^{(i)} = y^{(i)} (w^T x + b)$ , our fm margin will be doubled.

→ We can add normalization cond' that

$$\|w\|=1$$

### \* About Geometric Margin:



→ now the point  $x^{(i)} - \frac{\gamma^{(i)} w}{\|w\|}$  is on the hyperplane, it satisfies the eq<sup>n</sup> of the plane

so,

$$w^T \left( x^{(i)} - \frac{\gamma^{(i)} w}{\|w\|} \right) + b = 0$$

→ solve above eq<sup>n</sup> for  $\gamma^{(i)}$ ,

$$w^T x^{(i)} - \frac{\gamma^{(i)} w w^T}{\|w\|} + b = 0$$

$$w w^T = \|w\|^2, \text{ so}$$

$$\therefore w^T x^{(i)} + b = \frac{\gamma^{(i)} \|w\|}{\|w\|}$$

$$\therefore \gamma^{(i)} = \frac{(w)^T x^{(i)} + b}{\|w\|}$$

→ This signifies that if we've a training set  $(x^{(i)}, y^{(i)})$ , then the distance b/w  $x^{(i)}$  and separating hyperplane is defined by  $\gamma^{(i)}$ .

→ More generally, geometric margin:

$$\gamma^{(i)} = y^{(i)} \left[ \frac{\omega^T x^{(i)}}{\|\omega\|} + \frac{b}{\|\omega\|} \right]$$

→ We want geometric margin to be large.

→ If  $\|\omega\| = 1$ , then

$$\gamma^{(i)} = \hat{\gamma}^{(i)}$$

↑                      ↑  
functional          Geometric  
Margin               Margin.

→ So, more generally,

$$\boxed{\gamma^{(i)} - \frac{\hat{\gamma}^{(i)}}{\|\omega\|}}$$

→ This is correct understanding of both margins.

→ Geometric margin w.r.t entire training set is

$$\gamma = \min_i \gamma^{(i)}$$

Max. margin classifier :- It is the learning algo, which chooses the parameters  $w$  and  $b$  so to maximize geometric margin.

Best explanation of this SVN is on  
www.svn-tutorial.com.

Page No.	
Date	

$$\max_{\gamma, w, b} \gamma$$

such that

$$y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq \gamma, \|\mathbf{w}\|=1$$

this condition ~~is~~ ensures there is no training Geometric margin does not depend on  $\|\mathbf{w}\|$ .  
initially selected hyperplane margin does not change with data point b/w two points.

## Lecture-7

- Now the constraint  $\|\mathbf{w}\|=1$  is non-convex optimization constraint. We need to get rid of this constraint.
- So let's change the optimization problem.

since we know that

$$\gamma = \frac{\gamma}{\|\mathbf{w}\|}, \text{ we can write above problem as,}$$

$$\max_{\gamma, w, b} \frac{\gamma}{\|\mathbf{w}\|},$$

$$\text{subject to } y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq \gamma, i \in \{1, 2, \dots, m\}.$$

This is condition of constraint.

→ Now for simplicity, let's assume that

$$\gamma = 1,$$

so,

$$\min_i y^{(i)} (\omega^T x^{(i)} + b) = 1$$

→ Now rewriting the problem with  $\gamma = 1$

$$\min_{\omega, b} \frac{1}{2} \|\omega\|^2 \quad \begin{matrix} \text{max} \\ \text{same} \\ \text{as} \end{matrix} \quad \frac{1}{2} \|\omega\|^2$$

$$\text{subject to } y^{(i)} (\omega^T x^{(i)} + b) \geq 1$$

→ Now this is convex optimization prob with constraints. So we can use Lagrange multipliers to optimize.

### \* Lagrange Multipliers

$$\min_{\omega} f(\omega)$$

$$\text{subject to } h_i(\omega) = 0, i = 1, 2, \dots$$

Lagrangian:

$$L(\omega, \beta) = f(\omega) + \sum_i \beta_i h_i(\omega)$$

where  $\beta_i$  = Lagrangian Multipliers.

→ To solve it, set

$$\frac{\partial L}{\partial \omega} = 0 \quad \text{and} \quad \frac{\partial L}{\partial \beta} = 0$$

→ for some  $\omega^*$  to be a solution, it is necessary that

$\exists \beta^*$  such that

$$\frac{\partial L(\omega^*, \beta^*)}{\partial \omega} = 0 \quad \text{and}$$

$$\frac{\partial L(\omega^*, \beta^*)}{\partial \beta} = 0$$

→ In General,  $\xrightarrow{\text{primal problem}}$

$$\min_{\omega} f(\omega)$$

subject to  $g_i(\omega) \leq 0$ , for  $i=1, 2, \dots, k$

and  $h_i(\omega) = 0$  for  $i=1, 2, \dots, l$

Lagrangian:

$$L(\omega, \alpha, \beta) = f(\omega) + \sum_{i=1}^k \alpha_i g_i(\omega) + \sum_{i=1}^l \beta_i h_i(\omega) \quad \rightarrow ①$$

→ Define,

$$\text{Op}(\omega) = \max_{\alpha, \beta} L(\omega, \alpha, \beta)$$

$\alpha_i \geq 0$

Consider :

$$P^* = \min_{\omega} \max_{\alpha, \beta} L(\omega, \alpha, \beta)$$

$\alpha_i \geq 0$

$$= \min_{\omega} \text{Op}(\omega)$$

→ Primal Problem.

→  $\text{Op}(\omega) :=$  If  $g_i(\omega) > 0$ , then  $\text{Op}(\omega) = \infty$

how  $\text{Op}(\omega) = \infty$  ?

→ From previous page  $\sum_i \alpha_i g_i(\omega) \leq 0$ ,

$k$

$\sum_{i=1}^k \alpha_i g_i(\omega)$ , we can scale  $\alpha_i$  to arbitrarily large to make  $\text{Op}(\omega) = \infty$

→ And if  $h_i(\omega) \neq 0$ , then  $\text{Op}(\omega) = \infty$   
(Same reason)

→ otherwise,  $Q_p(w) = f(w)$  [By setting Lagrangian terms to zero].

∴ Thus,

$$Q_p(w) = \begin{cases} f(w) & \text{if constraints are satisfied } (g_i, h_i) \\ \infty & \text{otherwise} \end{cases}$$

So,  $\min_w Q_p(w) = \text{Original problem}$

### \* Dual Problem:

$$Q_D(\alpha, \beta) = \min_w L(w, \alpha, \beta)$$

$$d^* = \max_{\alpha \geq 0} \min_{w, \beta} L(w, \alpha, \beta) = \max_{\alpha \geq 0} Q_D(\alpha, \beta)$$

see the diff.  
with previous  
problem.

Dual  
problem.

$d^* \leq p^*$  is informally,

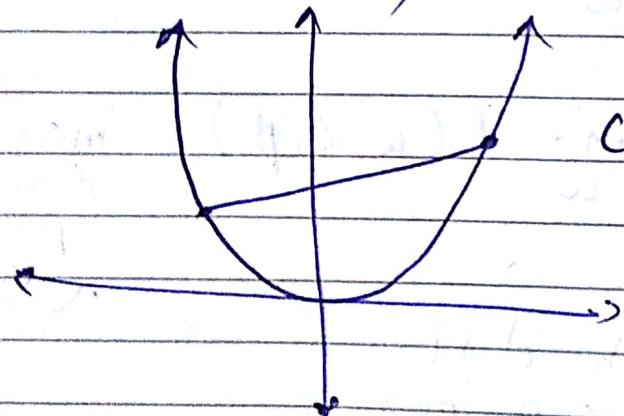
$$\max(\min(\dots)) \leq \min(\max(\dots))$$

- but sometime,

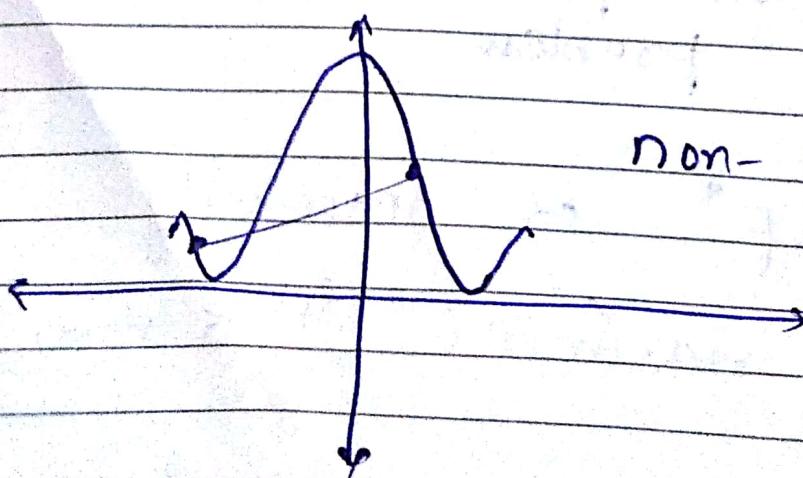
$$d^* = p^*$$

- It turns out sometimes it's easy to solve dual problem instead of primal problem.
- So here in SVM, we'll derive dual problem from primal problem and then solve dual problem.
- Conditions for which dual and primal problems give the same solution.

Let  $f$  be convex (Hessian  $H \geq 0$  or +ve semidefinite)



convex (Local minima  
is equal to global  
minima)



Non-Convex.

→ suppose,

$h_i$  is affine,  $[h_i(\omega) = \alpha_i^T \omega + b_i]$

and

$g_i$  is strictly feasible,  $[J\omega \text{ s.t } h_i g_i(\omega) < 0]$ .

→ Then  $\exists \omega^*, \alpha^*, \beta^*$  such that  $\omega^*$  solves the primal problem and  $\alpha^*, \beta^*$  solves dual problem, then

$$\beta^* = d^* = L(\omega^*, \alpha^*, \beta^*)$$

further,

~~$$\frac{\partial L}{\partial \omega} = 0 \quad \frac{\partial L}{\partial \alpha} = 0$$~~

$$\frac{\partial L}{\partial \beta} = 0$$

and  $\boxed{\alpha_i^* g_i(\omega) = 0}$  [KKT complementarity condition].

$$\hookrightarrow g_i(\omega^*) \leq 0 \text{ and } \alpha_i^* \geq 0$$

from this we can say that if product is zero, then one of the  $\alpha_i^*$  or  $g_i(\omega)$  is zero.

## \* Karush - Kuhn - Tucker complementarity (Cond) (KKT)

→ It implies that,

If  $\alpha_i^* > 0$  then  $g_i(\omega^*) = 0$

[bcz  $\alpha_i^* g_i(\omega^*) = 0$ ]

→ Usually,  $\alpha_i^* \neq 0 \Leftrightarrow g_i(\omega^*) = 0$  [Atleast one of them is zero].

→ When  $g_i(\omega^*) = 0$ , then,  $g_i(\omega)$  is an "active constraint".

\* Note:-

Lagrange Multipliers

Parameters

In  
Dual Problem  
 $\alpha_i, \beta_i$

In  
Actual S.V.  
 $\lambda_i$

$\omega$

$w, b$

→ Actual SVM minimization problem is,

$$\min_{\omega} \frac{1}{2} \|\omega\|^2 \quad [ \because \text{square is } ] \\ \text{s.t.} \quad [\text{for math. simplicity}] \\ y^{(i)} (\omega^T x^{(i)} + b) \geq 1, \quad i=1,2,\dots,m$$

→ so; constraint can be rewritten as,

$$g_i(\omega, b) = -y^{(i)} (\omega^T x^{(i)} + b) + 1 \leq 0$$

→ Implication of this constraint  $[g_i(\omega, b)]$  in terms of KKT conditions.

→ we've  $\alpha_i > 0 \Rightarrow g_i(\omega, b) = 0$  [Active const.]

$\Leftrightarrow (x^{(i)}, y^{(i)})$  has functional margin = 1.

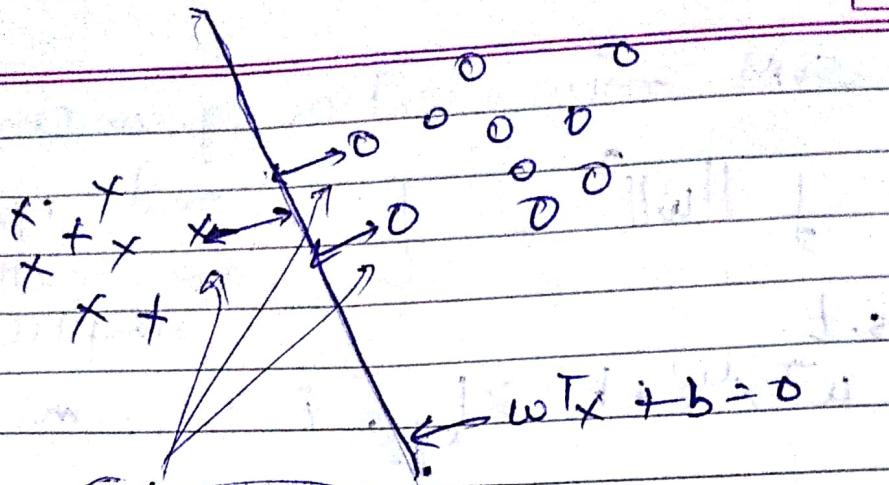
because;

$$g_i(\omega, b) = 0 = -y^{(i)} (\omega^T x^{(i)} + b) + 1$$

$$\therefore y^{(i)} (\omega^T x^{(i)} + b) = 1$$

∴ functional Margin = 1.

→ Geometrically what this means can be given by:-



functional margin = 1 and other examples have functional margin  $\geq 1$ .  
 Usually,  $\alpha_i \neq 0$  [ $\alpha_i > 0$ ]

→ Here the 3 points with functional margin = 1 are called support vectors. [few support vectors] so, usually  $\alpha_i = 0$  for most of the ~~support~~ vectors.

$$\rightarrow L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y^{(i)}(w^T x^{(i)}) + b) - 1$$

Dual problem:-

$$\mathcal{L}_D(\alpha) = \min_{w, b} L(w, b, \alpha)$$

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \underset{=} 0.$$

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\rightarrow \frac{\partial L}{\partial b} = - \sum_{i=1}^m \alpha_i y_i \stackrel{\text{set}}{=} 0$$

→ substitute those values back in ①

$$L = \frac{1}{2} w^T w - \sum_{i=1}^m \alpha_i (y^{(i)} (w^T x^{(i)} + b) - 1)$$

$$\left( \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T \left( \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right)$$

$$w^T w = \|w\|^2 = w \cdot w \quad [\text{Dot product}]$$

so,

$$L = \frac{1}{2} w \cdot w - \sum_{i=1}^m \alpha_i (y^{(i)} (w \cdot x + b) - 1)$$

$$= \frac{1}{2} \left( \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right) \left( \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right) -$$

$$\sum_{i=1}^m \alpha_i \left[ y_i \left( \left( \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right) \cdot x_i + b \right) - 1 \right]$$

$$= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)} x^{(j)} + \sum_{i=1}^m \alpha_i -$$

$$\sum_{i=1}^m \alpha_i y_i \left( \left( \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right) \cdot x_i + b \right)$$

$$\begin{aligned}
 &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i x_j - \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i x_j \\
 &\quad - b \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i \\
 &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j x_i x_j y_i y_j - b \sum_{i=1}^m \alpha_i
 \end{aligned}$$

Now,  $\sum_{i=1}^m \alpha_i y_i = 0$ ,

call it

$$w(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j x^{(i)} x^{(j)} y^{(i)} y^{(j)}$$

→ So, Dual problem is:

$$\max_w w(\alpha)$$

s.t

$$\alpha_i \geq 0$$

and  $\sum_{i=1}^m y_i \alpha_i = 0$

→ Interpretation of  $\sum_{i=1}^m y_i \alpha_i = 0$ ,

If  $\sum_i y_i \alpha_i \neq 0$ , then  $\Omega_D(\alpha) = -\infty$ , since we are maximizing  $\Omega_D(\alpha)$ , we better choose  $\sum_i y_i \alpha_i = 0$  to prevent penality.

→ So, if  $\sum_i y_i \alpha_i = 0$ , then  $\Omega_D(\alpha) = w(\alpha)$

→ solve  $w(x)$  for  $\alpha$ , then

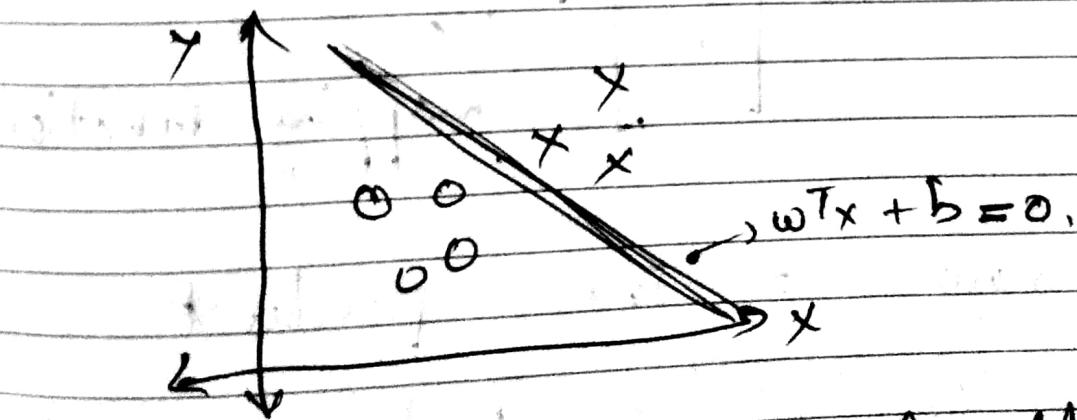
$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

→ then plug  $x$  and  $w$  to primal problem and solve for  $b$ ,

$$b = \max_{y_i=-1} (\omega^T x^{(i)}) + \min_{y_i=1} (\omega^T x^{(i)})$$

2.

→ Intuition is, Place the hyperplane in the exact middle of +ve & -ve example. [We've the orientation of plane given by  $w$ ,  $b$  just gives it's place]



→ This is called Optimal Margin Classifier or SVM.

## Lecture-8

Page No.

Date 29/5/18

→ As mentioned in dual problem,

$$\max_{\alpha} W(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

→ we only care about inner product or dot product of feature vectors  $x^{(i)}$  and  $\alpha^{(i)}$ , the actual values are not required.

### \* Idea behind Kernels:-

→ Let's say we've an attribute  $\alpha \in \mathbb{R}$  (for ex. living area of house).

→ Now take these feature and map it to a richer set of feature vector.

$\alpha \xrightarrow{\phi}$

$$\begin{bmatrix} x \\ x^2 \\ x^3 \\ x^4 \end{bmatrix}$$

$= \phi(x)$

mapping function.

→ So, now our inner product

$\langle x^{(i)}, \alpha^{(j)} \rangle$  will be replaced by

$$\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$

- Which is inner product in higher dimensions.
  - Some times  $\phi(x)$  is very high dimensional ( $\infty$  dimensional), so we can't compute this inner product efficiently.
  - Here Kernel ~~comes~~ comes into the picture
- $$K(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$$
- Kernels computes this inner product inexpensively, without actually transforming the original vector to higher dimensions. [without mapping]

\* How this is done using Kernels?

- Let's say we've two vectors  $x, z \in \mathbb{R}^n$  and ~~kernel~~ kernel defined as,

$$K(x, z) = (x^T z)^2 = \left( \sum_{i=1}^n x_i z_i \right) \left( \sum_{j=1}^n x_j z_j \right)$$

$$= \sum_{i=1}^n \sum_{j=1}^n (x_i x_j) (z_i z_j)$$

$$= (\phi(x))^T (\phi(z))$$

→ where,

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

$\phi(z) = \text{same as } \phi(x)$ .

for  $n=3$ .

→ So, this  $(x^T z)^2$  kernel calculates the inner product without even transforming or mapping the original vector to  $\phi$ . So it is highly efficient idea.

→ So, we need  $O(n^2)$  to compute just  $\phi(x)$ .

→ And to compute  $k(x, z)$ , it needs  $O(n)$  time.

→ One general example of kernels is,

$K(x, z) = (x^T z + c)^d$ , which is  $\binom{n+d}{d}$  features of all monomials upto degree  $d$ .

- \* How to come up with a good kernel in actual machine learning problems?
- if we've two feature vector and mappings,  
 $x \rightarrow \phi(x)$ ,  $z \rightarrow \phi(z)$ , and  
 $\langle \phi(x), \phi(z) \rangle$  is inner product in ~~in~~ higher dimensions, then one intuition is,
- if  $x$  and  $z$  are similar kind of features (vectors), then they are pointing in similar direction, so inner product or dot product will be higher because of

$$x \cdot z = |x| |z| \cos \theta \text{ and}$$

$$\cos \theta \rightarrow 1, \text{ when } \theta \rightarrow 0.$$

- And similarly, if data or vectors are very dissimilar, then dot product will be very low.

So, for particular problem, choose,

$$k(x, z) \begin{cases} \text{Large if } x, z \text{ are similar} \\ \text{Small, if } x, z \text{ are dissimilar.} \end{cases}$$

Gaussian  
Kernel

$$\text{★ } K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right), \text{ so,}$$

$\exists \phi$  such that  $K(x, z) = \langle \phi(x), \phi(z) \rangle$   
or in other words, if this kernel  
valid?

→ Suppose  $K$  is kernel., then let

$\{x^{(i)}, \dots, x^{(m)}\}$  be given. and,

Let  $K \in \mathbb{R}^{m \times m}$  [K is matrix],  
such that

$K_{ij} = K(x^{(i)}, x^{(j)})$ , then for  
any vector  $z \in \mathbb{R}^m$ ,

consider,

$$z^T K z = \sum_i \sum_j z_i K_{ij} z_j$$

$$= \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(j)}) z_j$$

$$= \sum_i \sum_j z_i \sum_k (\phi(x^{(i)}))_k (\phi(x^{(j)}))_k$$

$$= \sum_k \sum_i \sum_j z_i (\phi(x^{(i)}))_k (\phi(x^{(j)}))_k z_j$$

$$= \sum_{i,j} K_{ij} (\sum_i z_i \phi(x^{(i)})^T \phi(x^{(j)}))^2 \geq 0.$$

So,  $K \geq 0$ , so matrix  $K$  is positive semi-definite matrix.

→ So, from this, we can say that if  $K$  is valid kernel, then matrix  $K$  is tve semi-definite.

→ the converse also holds true: that is if matrix  $K$  is tve semi-definite then the kernel corresponding to it is valid kernel.

→ So, Theorem (Mercer): Let  $K(x, z)$  be given, then  $K$  is valid kernel, or mercer kernel (i.e.  $\exists \phi$  s.t  $K(x, z) = \phi(x)^T \phi(z)$ ),

iff for any set of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  ( $m < \infty$ ), the kernel matrix  $K \in \mathbb{R}^{m \times m}$  is symmetric tve semidefinite.

→ Converse is also true as mentioned previously.

→ Ex. of invalid kernel,

$$K(x, z) = -1 \neq \phi(x)^T \phi(z)$$

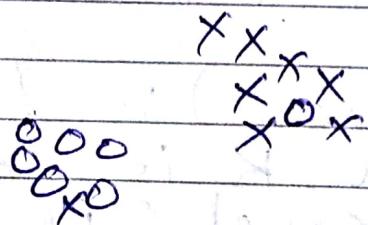
because inner product always  $\geq 0$ .

→ Intuition for this kernel is given best in SVM book.

### \* L1 Norm Soft Margin SVM :-

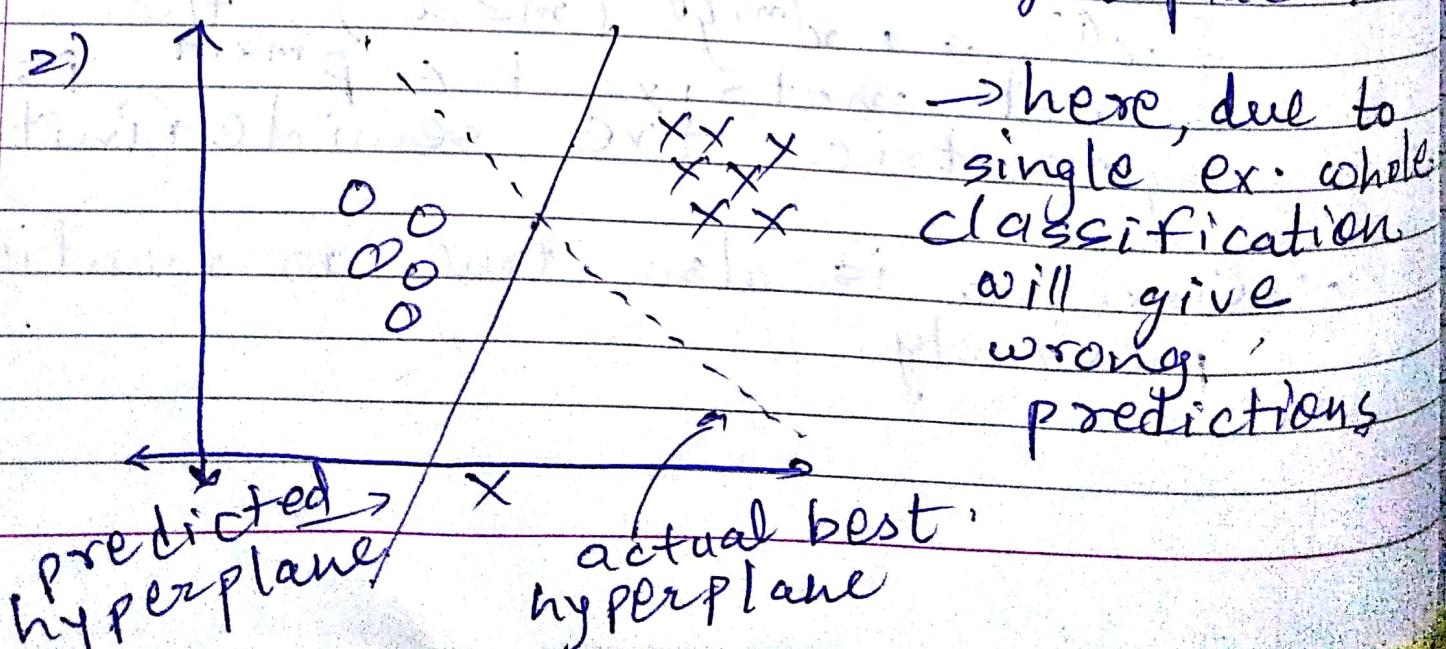
→ There are few cases, where hard margin SVM does not work well.

1)



⇒ Here hard margin SVM won't figure out the hyperplane, even if only one ex. is not at right place.

2)



→ here, due to single ex. whole classification will give wrong predictions

↓

actual best  
hyperplane

$\zeta$  (zeta).

Page No.	
Date	

→ The updated formulations for soft margin SVM is,

$$\min_{w, b, \xi_i} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

such that

functional margin  $\begin{cases} y^{(i)} (w^T x^{(i)} + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \text{ for any } i = 1, 2, \dots, m \end{cases}$

→ If  $y^{(i)} (w^T x^{(i)} + b) > 0$  that means SVM has classified the example correctly. and if it's  $\leq 0$ , we miss classified it. so, here we can set some of  $\xi_i \geq 1$ , we are allowing SVM to miss classify for some of the examples.

→ However, it is not encouraged, as we are adding penalty to our final goal by

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

from this, the more  $\xi_i \geq 1$ , less minimized our goal will be.

→ Here 'C' controls the effect of penalty.

→ Now, corresponding to the set of new constraints, our new generalized Lagrangian will be,

$$L(\omega, b, \xi, \alpha, \gamma) = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m \xi_i$$

$$- \sum_{i=1}^m \alpha_i (y^{(i)} (\omega^T x^{(i)} + b) - 1 + \xi_i) - \sum_{i=1}^m \gamma_i \xi_i$$

[from  $L(\omega, \alpha, \beta) = f(\omega) + \sum_{j=1}^k \alpha_j g_j(\omega) + \sum_{j=1}^l \beta_j h_j(\omega)$ ]

→ so, here  $\alpha$  and  $\gamma$  are Lagrange multipliers, and

$$g_i(\omega) = y^{(i)} (\omega^T x^{(i)} + b) - 1 + \xi_i \quad \left. \right\} \text{constraint}$$

$$h_i(\omega) = \sum_{j=1}^m \xi_j$$

$$\text{and } f(\omega) = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m \xi_i$$

→ And from this we can derive dual of this optimization problem same as previously, we get following.  $f^*$ .

$$\max_{\alpha} w(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

~~subject to~~

subject to

$$\sum_{i=1}^m y^{(i)} \alpha_i = 0 \text{ and}$$

$$0 \leq \alpha_i \leq C \text{ for } i = 1, 2, \dots, m.$$

Note that only  
this changes compared to  
previous dual problem.

→ convergence conditions,

$$\alpha_i = 0 \Rightarrow y^{(i)} (\omega^T x^{(i)} + b) \geq 1$$

$$\alpha_i = C \Rightarrow y^{(i)} (\omega^T x^{(i)} + b) \leq 1$$

$$0 < \alpha_i < C \Rightarrow y^{(i)} (\omega^T x^{(i)} + b) = 1$$

→ To actually solve the above optimization problem, we need efficient algo (CSMO).

\* Digression :- Coordinate ascent :-

$\max \omega(\alpha_1, \alpha_2, \dots, \alpha_m)$  [No constraints on  $\alpha_i$ 's]

- In coordinate ascent, maximization is done by maximizing  $\omega$  along one direction at a time or using one  $\alpha_i$  at a time.  
I.e. solving univariate opt. problem in loop J.

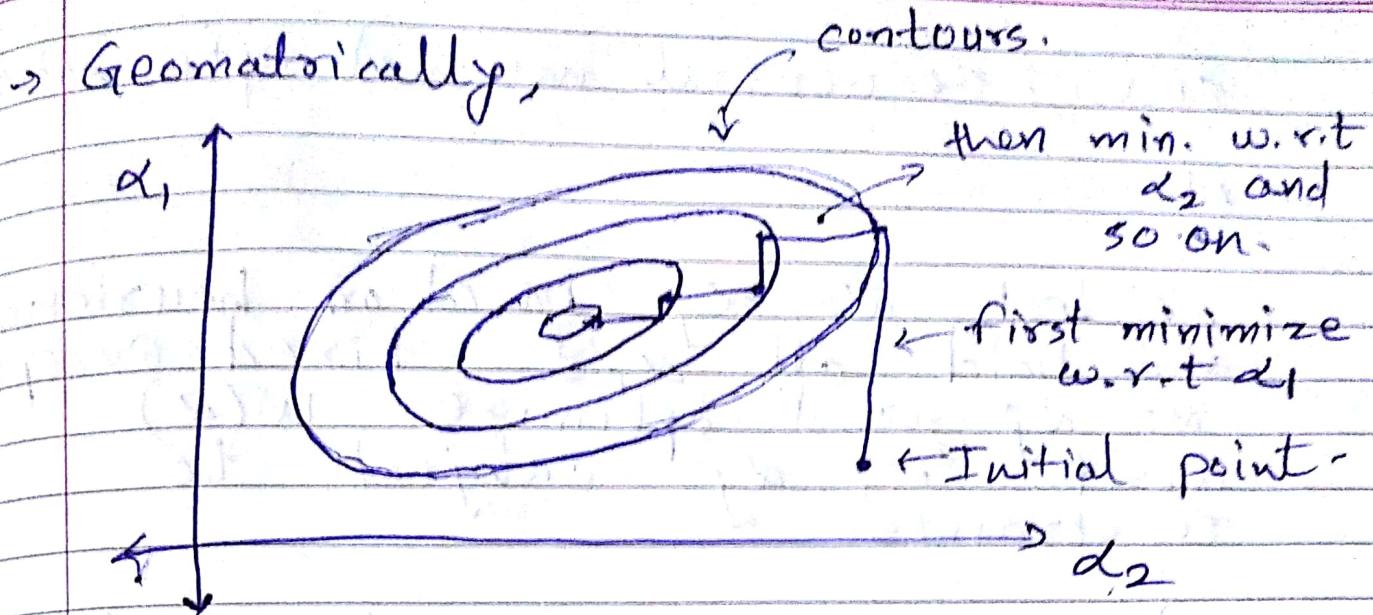
- Algo. can be given by :-

Repeat {

for  $i=1$  to  $m$

$$\alpha_i^* = \underset{\alpha_i}{\operatorname{argmax}} \omega(\alpha_1, \dots, \hat{\alpha}_i, \dots, \alpha_m)$$

Informally this means hold everything except  $\alpha_i$  fixed.



→ Compared to Newton's method, this take lot more steps, but the main advantage is sometime the object  $W(\alpha_i)$  is very inexpensive w.r.t any one of the parameters ( $\alpha_i$ ).

→ Now, for applying this algo. to our SVM softmargin problem, the problem is ~~that~~ with basic form of this algo., having no constraints over  $\alpha_i$ . But we've constraints over  $\alpha_i$  in softmargin SVM., which is

$$\sum_{i=1}^m y^{(i)} \alpha_i \geq 0.$$

→ So, changing 1  $\alpha_i$  at a time won't work. Now we'll change 2  $\alpha_i$  at a time, which algo. is called

## SMO (Sequential minimal optimization)

→ Algo:-

Select  $\alpha_i, \alpha_j$  based on heuristic and hold all  $\alpha_i$ 's fixed except  $\alpha_i, \alpha_j$  and optimize  $W(\alpha)$  w.r.t.  $\alpha_i, \alpha_j$  subject to constraints.

→ And keep running the algo. until the convergence criteria is met. [convergence cond' as mentioned previously].

\* How  $\alpha_i$  and  $\alpha_j$  are updated to optimize  $W(\alpha)$ ?

→ In every step, we respect the constraint

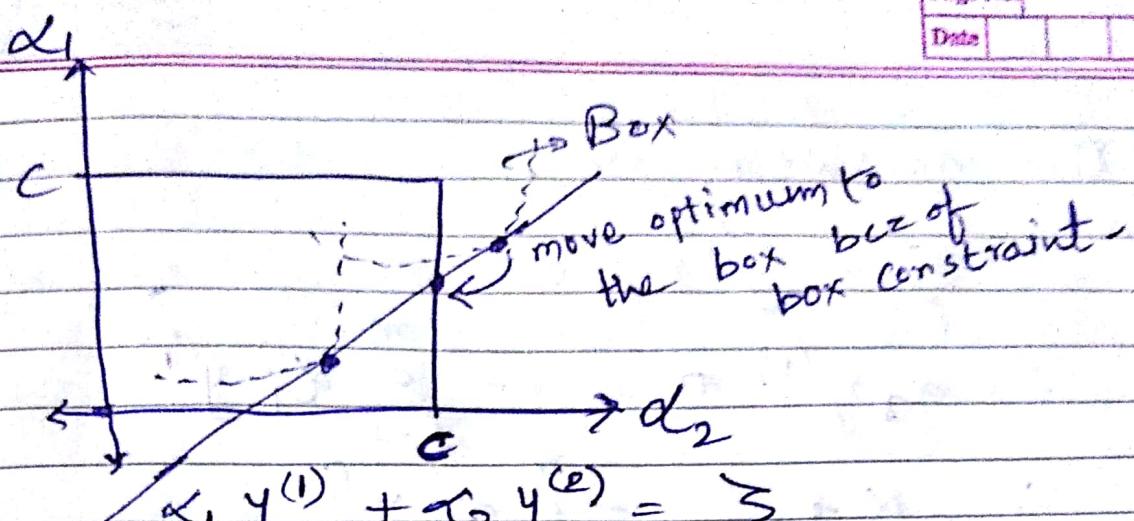
$$\sum_i \alpha_i y^{(i)} = 0$$

→ let's consider  $i=1$  and  $j=2$  for sake of simplicity, so,

$$\alpha_1 y_1 + \alpha_2 y_2 = - \sum_{i=3}^m \alpha_i y^{(i)} = \zeta$$

~~other~~

→ also other constraint is  $0 \leq \alpha_i \leq C$ . [Box constraint]



$$\text{so, } \alpha_1 = \frac{3 - x_2 y^{(2)}}{y^{(1)}}$$

$$w(\alpha_1, \alpha_2, \dots, \alpha_m)$$

$$= w\left(\frac{3 - x_2 y^{(2)}}{y^{(1)}}, \alpha_2, \alpha_3, \dots, \alpha_m\right)$$

→ Now we've  $\alpha_3, \alpha_4, \dots, \alpha_m$  are fixed,  
so this  $f^n$  is only 1-D  
quadratic function of  $\alpha_2$ .

$$= a\alpha_2^2 + b\alpha_2 + c$$

→ And by optimizing above  $f^n$ ,  
will give optimized value of  
 $w(\alpha_1, \alpha_2, \dots, \alpha_m)$  w.r.t  $\alpha_2$ .

→ As underlined on previous page,  
why changing one  $\alpha_i$  at a  
time won't work for  
constrained dis?

→ The answer is,

$$\begin{aligned} \sum_i \alpha_i y^{(i)} &= 0 \\ \therefore \alpha_1 y^{(1)} &= - \sum_{i=2}^m \alpha_i y^{(i)} \\ \therefore \alpha_1 &= \frac{- \sum_{i=2}^m \alpha_i y^{(i)}}{y^{(1)}} \end{aligned}$$

→ So,  $\alpha_1$  is written as function of  $\alpha_2, \alpha_3 \dots \alpha_m$  and if we keep  $\alpha_2, \alpha_3 \dots \alpha_m$  fixed, then we can't change  $\alpha_1$ .

→ But in case of changing 2  $\alpha_i$  ( $\alpha_1$  and  $\alpha_2$  for example), then if we change  $\alpha_1$ ,  $\alpha_2$  will adjust accordingly without affecting other  $\alpha_3, \alpha_4 \dots \alpha_m$ .

→ And this is the reason why we have to change 2  $\alpha_i$  at a time in constrained problem.

→ We are changing minimum "2"  $\alpha_i$ 's at a time and that's why it is called sequential minimization.

## → Applications of SVMs -

→ In handwritten digit recognition, it gives same performance as best trained Neural network.



→ This results best when choosing kernel as following,

$$K(x, y) = (x^T y)^d \quad \text{or}$$

$$\text{gaussian kernel } K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$$

→ Another application is in classification of proteins.

→ Protein sequence is sequence of Amino Acids, represented by A, B, C, ..., Z.

→ So, for a random length protein sequence,

BAJTSIAJBASTAU ..

→ our feature vector  $\phi(x)$  will be constructed by writing all comb' of length 4 amino acids.

$$\begin{array}{c}
 \text{AAAA} \\
 \text{AAAB} \\
 \vdots \\
 \text{BAJT} \\
 \vdots \\
 \text{TSTA} \\
 \vdots \\
 \text{ZZZZ}
 \end{array}
 \left| \begin{array}{c} 0 \\ 0 \\ \vdots \\ 2 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{array} \right| \cdots \cdots \cdots = \phi(x).$$

$$\text{So, } \phi(x) \in \mathbb{R}^{(2^4)} \quad 160000$$

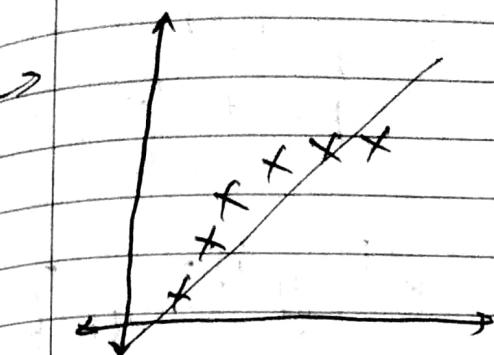
→ Very high dimensional feature vector. but can be calculated  $\phi(x)^T \phi(x)$  by efficient dynamic programming algo.

# Lecture - 9

Page No.  
Date 3/5/18

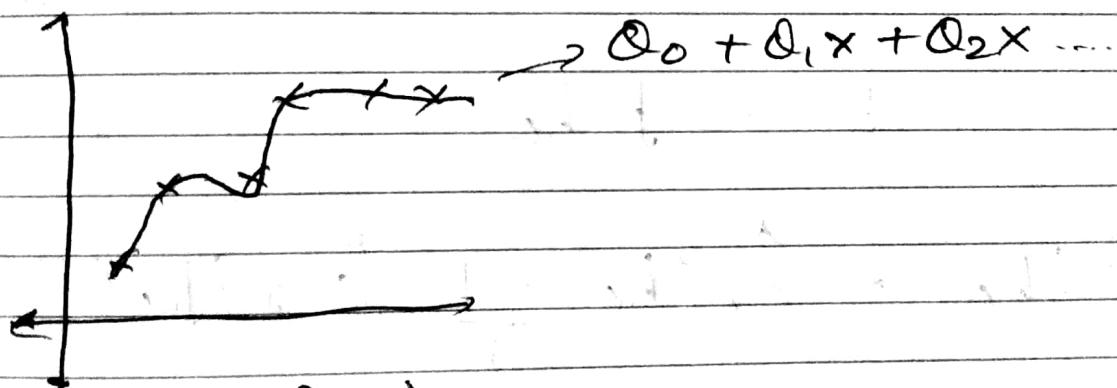
## \* Learning Theory:-

→ Bias - Variency trade off.

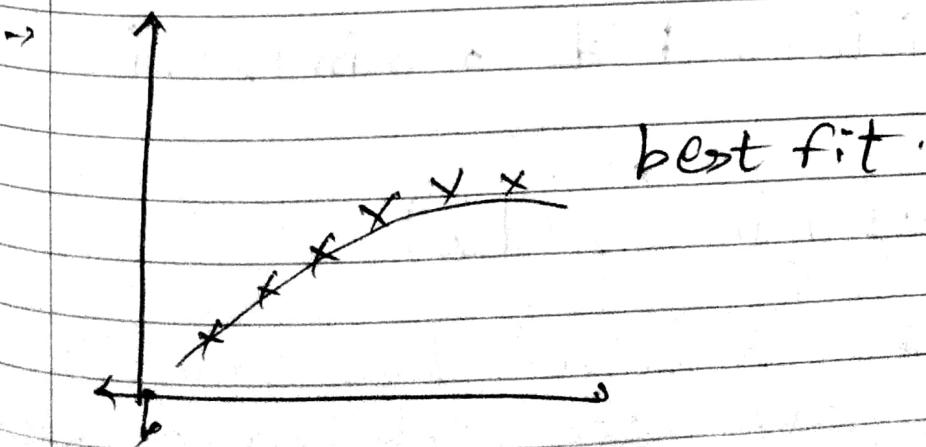


→ This fails to fit the quadratic nature of the data so it is unfit or bias.

$Q_0 + Q_1x$   
"unfit" or  
"bias"



overfitting  
or high variance



## \* Linear classifier:

$$h_0(x) = g(\theta^T x)$$

$$g(z) = \underbrace{1}_{\text{nothing but}} \{z \geq 0\}, \quad (y \in \{0, 1\})$$

perceptron learning algo.

$$S = \{(x^{(i)}, y^{(i)})\}_{i=1}^m \quad (x^{(i)}, y^{(i)}) \sim \text{IID}$$

training set

some distribution like gaussian

→ Training Error of  $h_0$ :

$$\hat{\epsilon}(h_0) = \hat{\epsilon}_s(h_0) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h_0(x^{(i)}) \neq y^{(i)}\}$$

→ Training Error is also called risk.

\* ERM (Empirical Risk minimization):

$$\hat{\theta} = \arg \min_{\theta} \hat{\epsilon}_s(h_0)$$

→ This problem is solved by logistic regression, SVM etc. learning algo.

→ This is non-convex optimization problem.  
 But SVM, Logistic regression solved this ~~by~~ as convex optimization problem.

\* Hypothesis class  $H = \{h_\theta : \theta \in \mathbb{R}^{n+1}\}$ .

→ class of all hypothesis, from which all learning algo. choose from.

Each  $h_\theta : X \rightarrow \{0, 1\}$

$\xrightarrow{\text{input domain}}$	$\xrightarrow{\text{range}}$	}
→ class of all linear classifiers	→ In general it can be any set like class of neural networks	

→ So, ERM can be written,

$$\hat{h} = \arg \min_{h \in H} \hat{E}_S(h)$$

→ In training a model, our goal is not to minimize error while predicting on training data set. But our goal is to make correct predictions (min. error) on test data set (the data which is unseen to the model while training).

→ So, Generalization Error :-

$$E(h) = P_{(x,y) \sim D} (h(x) \neq y)$$

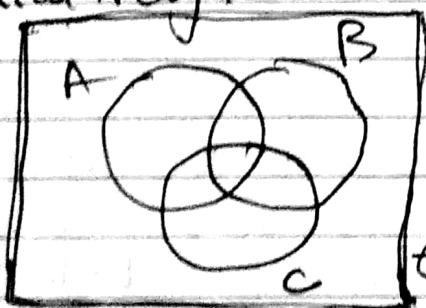
new example      ↗ mistable  
the new example.

→ Some imp. lemma:-

i) Union Bound :- Let  $A_1, A_2, A_3, \dots, A_k$  be  $k$  events (not necessarily independent), then

$$P(A_1 \cup A_2 \cup A_3 \dots \cup A_k) \leq P(A_1) + P(A_2) + \dots + P(A_k)$$

→ Intuitively :-



→ Area of the ~~whole~~ whole shape is less than equal to area of the individual circles A, B, C.

that on any symbol is estimation of  
that. for ex.  $\hat{E}, \hat{\phi}$  etc.

Page No.		
Date		

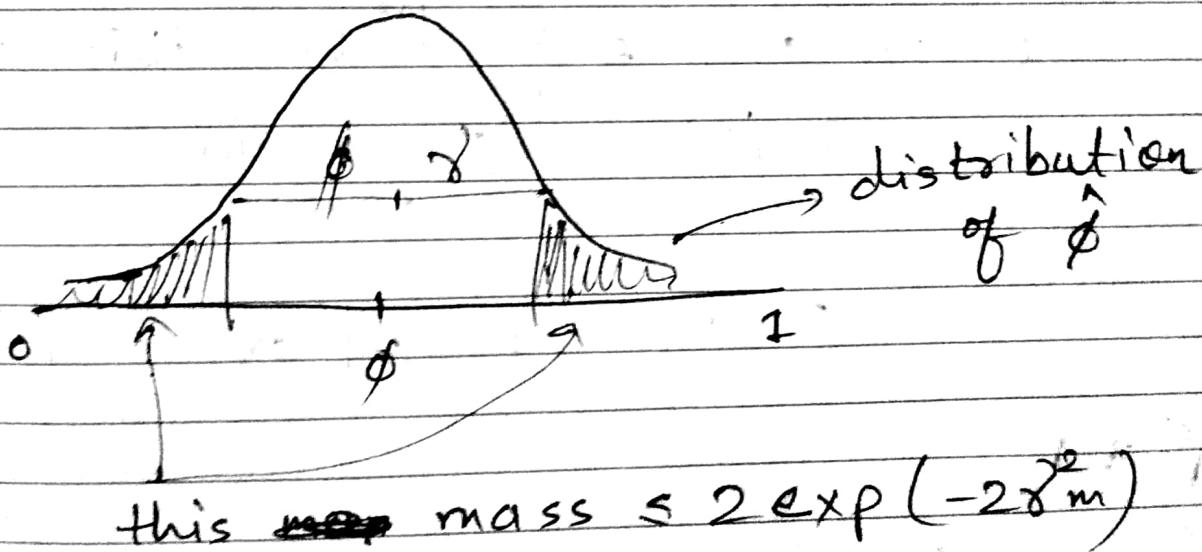
2) Hoeffding Inequality: Let  $z_1, z_2, \dots, z_m$  be  $m$  IID Bernoulli ( $\phi$ ) random variables ( $P(z_i=1)=\phi$ ),

Let  $\hat{\phi} = \frac{1}{m} \sum_{i=1}^m z_i$  and let for

any  $\gamma > 0$  be fixed, then

$$P(|\hat{\phi} - \phi| > \gamma) \leq 2 \exp(-2\gamma^2 m)$$

→ Intuitively:



→  $2 \exp(-2\gamma^2 m)$  decreases exponentially with  $m$ . bcz  $\gamma$  is fixed.

→ ERM for finite hypothesis class  $H \subseteq$

$H = \{h_1, h_2, \dots, h_k\}$ ,  $k$  hypotheses.

Each  $h_i$  is a function

$$h = \arg \min_{h_i \in H} \hat{\epsilon}_s(h_i)$$

↳ This means ERM picks the  $h_i$  with lowest training error.

→ Need to prove that minimizing training error also minimizes the generalization errors.

Strategy:- 1)  $\hat{\epsilon} \approx \epsilon$  <sup>first</sup> show this result.  
 ↓      ↑  
 training    generalizing  
 error        error.

2) Show bound on  $\hat{\epsilon}(h)$

<sup>also</sup>  
 generalization  
 errors

~~1st step~~

→ So, According to above steps, let's start.  
 fix proof.

→ fix any  $h_j \in H$ .

Define  $z_i = \mathbb{I}\{h_j(x^{(i)}) \neq y^{(i)}\} \in \{0, 1\}$

Training error is also called in-sample error, and generalization is also called out-sample error.

Date		
------	--	--

$$P(z_i=1) = E(h_j) \rightarrow \begin{array}{l} \text{generalization} \\ \text{error in } h_j. \end{array}$$

probability  
of misclassification.

→ so,  $z_i$  is bernoulli random variable with mean  $E(h_j)$ .

→ why mean is  $E(h_j)$ ?

bcz according to bernoulli distribution,

$$P(x=1) = p \text{ and}$$

$$P(x=0) = 1-p.$$

so, mean = expected value.

$$\text{so, } E[X] = P(X=1) \cdot 1 + P(X=0) \cdot 0 \\ = p \cdot 1 + q \cdot 0$$

$$E[X] = p = \text{mean.}$$

$$\rightarrow \text{training error } \hat{E}(h_j) = \frac{1}{m} \sum_{i=1}^m z_i$$

$$= \frac{1}{m} \sum_{i=1}^m \sum_{j \in \{h_j(x^{(i)}) \neq y^{(i)}\}}$$

$$\rightarrow \text{mean of each } z_i = E(h_j) \\ = \text{generalization error.}$$

→ By Hoeffding's inequality:-

$$P(|\hat{\epsilon}(h_j) - \hat{\epsilon}(h_j)| > \delta) \leq 2 \exp(-2\delta^2 m)$$

→ Intuitively this means that when we've large training examples, then the  $m$  will be large and the R.H.S will be small.

→ So, the absolute difference between training error and generalization error is very small.

$$\text{So, } P(|\hat{\epsilon}(h_j) - \hat{\epsilon}(h_j)| > \delta) \approx 0$$

→ We've proved this for a fix hypothesis  $h_j$ , now we want to prove it for entire  $H$  or for  $h_i$ ,  $i = 1, 2, \dots, k$ .

→ Let define  $A_j$  = event such that

$$|\hat{\epsilon}(h_j) - \hat{\epsilon}(h_j)| > \delta$$

$$\text{So, } P(A_j) \leq 2 \exp(-2\delta^2 m)$$

$$\rightarrow P(\exists h_j \in H, |\hat{\epsilon}(h_j) - \hat{\epsilon}(h_j)| > \delta)$$

$$= P(A_1 \cup A_2 \cup \dots \cup A_k) \leq \sum_{i=1}^k P(A_i)$$

[By Union Bound]

$$\leq \sum_{i=1}^k 2 \exp(-2\delta^2 m)$$

and  ~~$\sum$~~   $\sum_{i=1}^k 2 \exp(-2\delta^2 m)$

$$= 2 \exp(-2\delta^2 m) \sum_{i=1}^k 1.$$

$$= 2k \exp(-2\delta^2 m)$$

→ Taking negation  $(1 - P)$  at both sides,

$$1 - P(\exists h_j \in H / |\hat{e}(h_j) - \hat{e}(h_j)| > \gamma)$$

$$\geq 1 - 2k \exp(-2\delta^2 m)$$

$$\therefore P(\text{not } \exists h_j \in H / |\hat{e}(h_j) - \hat{e}(h_j)| > \gamma)$$

$$\geq 1 - 2e^{-2\delta^2 m}$$

$$\therefore P(\forall h_j \in H / |\hat{e}(h_j) - \hat{e}(h_j)| \leq \gamma)$$

$$\geq 1 - 2ke^{-2\delta^2 m}$$

→ So, with probability  $1 - 2ke^{-2\delta^2 m}$ ,  $\hat{e}(h)$  will be within  $\gamma$  of  $e(h)$  for all  $h \in H$ .

→ This is called Uniform convergence.  
 →  $\gamma$  is constant and above result is true for any value of  $\gamma$ .

$m$  = training set size

Page No.	
Date	

- here we proved that given  $\gamma$  and  $m$ , what is the probability of uniform convergence.
- Another equivalent form is given  $\gamma$  and probability  $\delta$  (delta) of making a large error, what should be the value of  $m$ ?

$$\delta = 2k e^{-2\gamma^2 m}$$

$$\therefore \frac{\delta}{2k} = e^{-2\gamma^2 m}$$

$$\therefore \frac{2k}{\delta} = e^{2\gamma^2 m}$$

$$2\gamma^2 m = \ln \frac{2k}{\delta}$$

$$\therefore m = \frac{1}{2\gamma^2} \ln \frac{2k}{\delta}$$

So long as  $m \geq \frac{1}{2\gamma^2} \ln \frac{2k}{\delta}$ , then

with probability  $1 - \delta$ , we've

$$|E(h) - \hat{E}(h)| \leq \gamma \text{ for all } h \in H.$$

- This is called "Sample Complexity" bound.

→ In practice

$$\forall k, \log k \leq 30$$

\* Another form is "Error bound", here we solve for  $\gamma$ .

$$\therefore \gamma = \sqrt{\frac{L}{2m} \ln \frac{2k}{\delta}}$$

→ So with probability 1- $\delta$ , we've that  $\forall h \in H$ ,

$$|\epsilon(h) - \hat{\epsilon}(h)| \leq \sqrt{\frac{L}{2m} \ln \frac{2k}{\delta}} = \gamma$$

→ step

Let's assume  $\forall h \in H, |\epsilon(h) - \hat{\epsilon}(h)| \leq \gamma \rightarrow (1)$

→ Can we prove something about  $\hat{\epsilon}(h)$ ? where  $\hat{h} = \operatorname{argmin}_{h \in H} \hat{\epsilon}(h) \rightarrow (2)$

→ Also let's define  $h^* = \operatorname{argmin}_{h \in H} \epsilon(h) \rightarrow (3)$

so  $h^*$  is best possible hypothesis in case of min. generalization error

$$\begin{aligned}
 E(\hat{h}) &\leq \hat{E}(\hat{h}) + \gamma \quad \text{by (1)} \\
 &\leq \hat{E}(h^*) + \gamma \quad \text{by (2)} \\
 &\leq (E(h^*) + \gamma) + \gamma \\
 &\leq E(h^*) + 2\gamma \quad \text{by (1)}
 \end{aligned}$$

→ here  $\hat{h}$  is the hypothesis chosen by ERM. [i.e. the hypothesis, which gives min. error].

- $\hat{E}(h)$  = training error of  $h$
- $E(h)$  = generalization error of  $h$ .

\* Theorem: Let given finite set of  $k$  hypotheses  $|H| = k$ , and let  $m, s$  be fixed, then with probability of  $1-s$ ,

$$E(\hat{h}) \leq \underbrace{\left( \min_{h \in H} E(h) \right)}_{E(h^*)} + 2 \sqrt{\frac{1}{2m} \log \frac{2k}{s}}$$

Set  $\gamma = \sqrt{\frac{1}{2m} \log \frac{2k}{s}}$ , we know

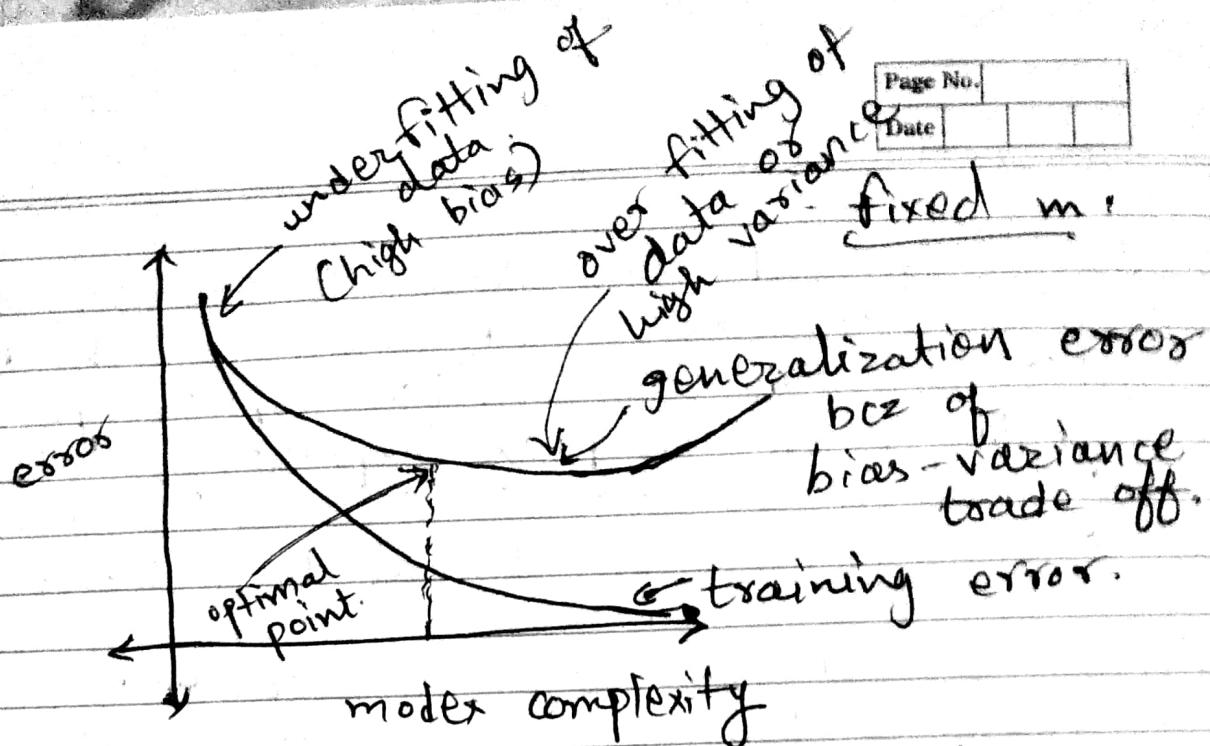
$|E(h) - \hat{E}(h)| \leq \gamma$  for  $\forall h \in H$   
 holds with probability  $1-s$  which  
 implies \*

- we've proved this for linear hypothesis class  $H$ . but if we switch from linear to quadratic hypothesis (to more complex hypothesis class), generalization error  $E(h^*)$  will always be less than or equal to linear class.
- So, switching to more complex class decreases the error in generalization
- But  $k$  increases. so generalization error inc. with  $k$ , as shown in eqn (\*)
- This is called bias-variance trade off. It is in the hope of choosing min. generalizing error we tend to choose more complex hypothesis class, but doing so increases value of  $k$ , which intern increases the gen. error

$$E(h^*) \leq \left( \min_{h \in H} E(h) \right) + 2 \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}$$

"bias"                          "variance"

→ so most complex class  $\rightarrow$  bias  $\downarrow$   
 variance  $\uparrow$



[Degree of Polynomial,  
size of  $H$  etc]

\* Corollary :- Let  $|H|=k$ , let any  $\delta, \gamma$  be fixed, Then for

$$E(\hat{h}) \leq \min_{h \in H} E(h) + 2\gamma$$

with probability  $1-\delta$ , if suffices that

$$m \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta} = O\left(\frac{1}{\gamma^2} \log \frac{k}{\delta}\right)$$

## Lecture 10

Page No.  
Date 3 6 18

→ Now we need to generalize the previous corollary for hypothesis class.

→ Let's say  $H$  is parameterised by  $d$  real numbers.

[So, if we are applying logistic regression with  $n$  features, then  $d = n+1$ ]

→ Computer uses 0 and 1 bits to represent real numbers. And each real no. is represented by .. 64 bit representation.

→ So  $H$  is parameterised by  $64d$  bits.

$$|H| = 2^{64d} \quad [\text{no. of ways to flip } 64d \text{ bits}]$$

→ Suffices that

$$m \geq O\left(\frac{1}{\gamma^2} \log \frac{2^{64d}}{\delta}\right)$$

$$\geq O\left(\frac{d}{\gamma^2} \log \frac{1}{\delta}\right)$$

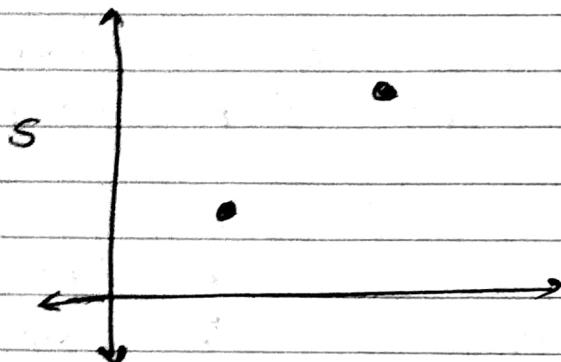
→ Definition: Given a set  $S = \{x^{(1)}, x^{(2)}, \dots, x^{(d)}\}$ , we say  $H$  shatters  $S$ , if  $H$  can realize any ~~lets~~ labeling on it.

→ "H shatters S" means H has some hypothesis  $h \in H$ , which labels that 'd' examples

correctly.

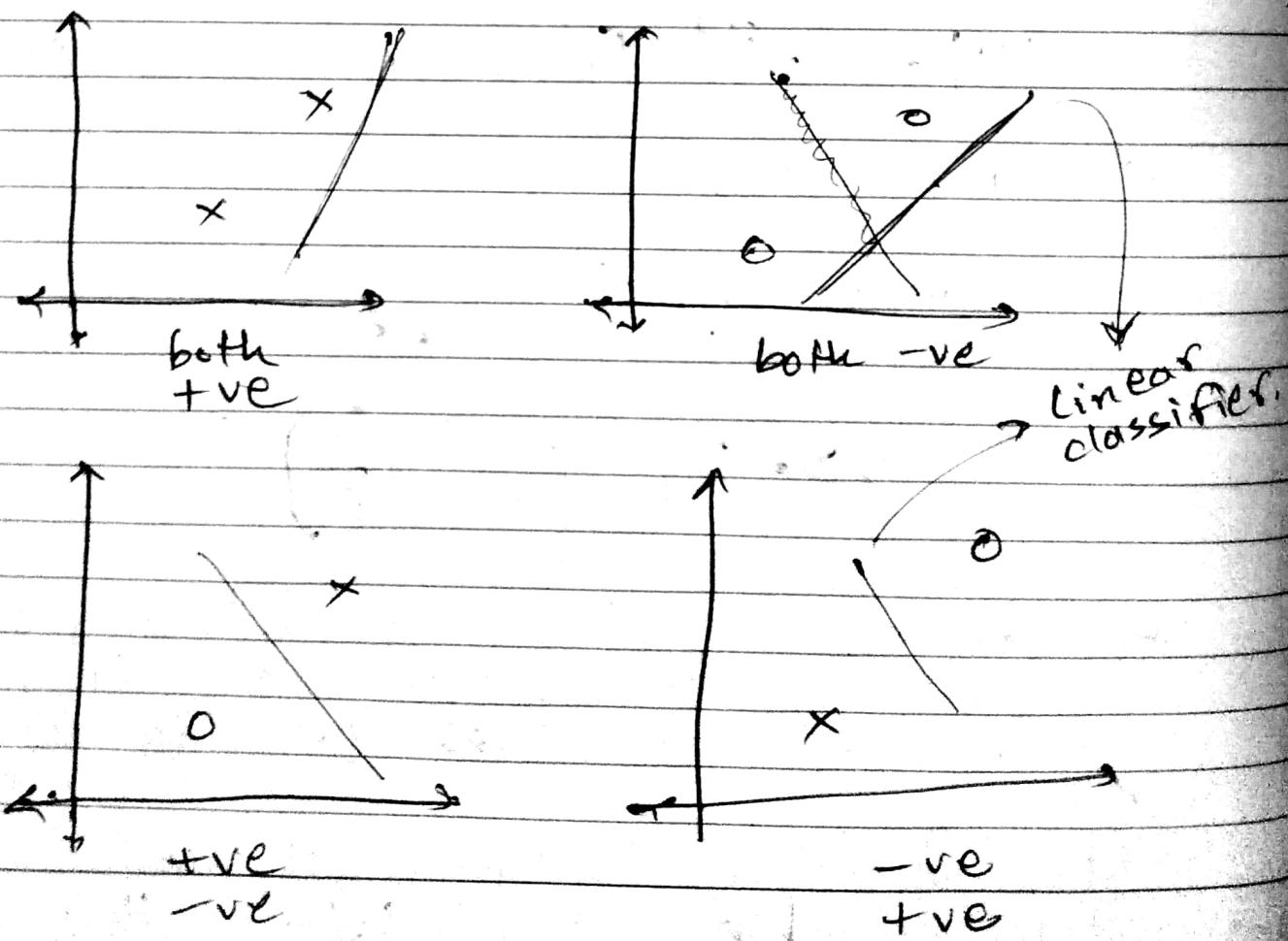
→ For Example :-

$H = \{ \text{linear classifiers in } 2D \}$



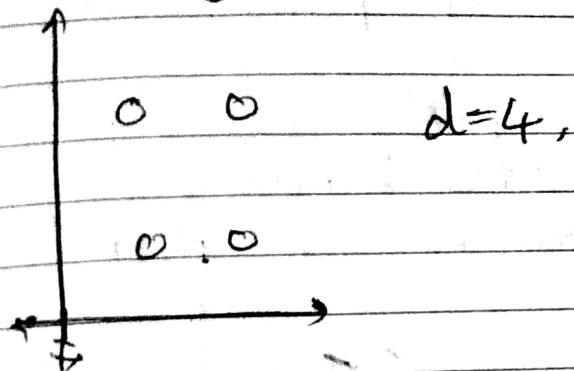
→ Set  $s$  contains  
2 points as  
shown.

→ Then there are 4 possible ways to  
label them.

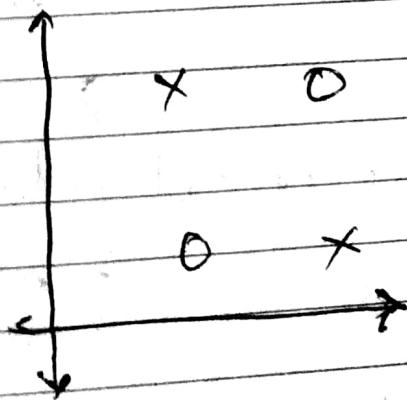


→ So, this shows H class shatters the set S of 2 points. ( $d=2$ )

→ So for another example, if we have ~~set~~ following S,



→ Then H of linear hypothesis does not shatter S, because there is no linear boundary for some sets.



\* Definition: The Vapnik and Chervonenkis dimension (VC dimension) of H (VCH) is the size of largest set S, which is shattered by H.

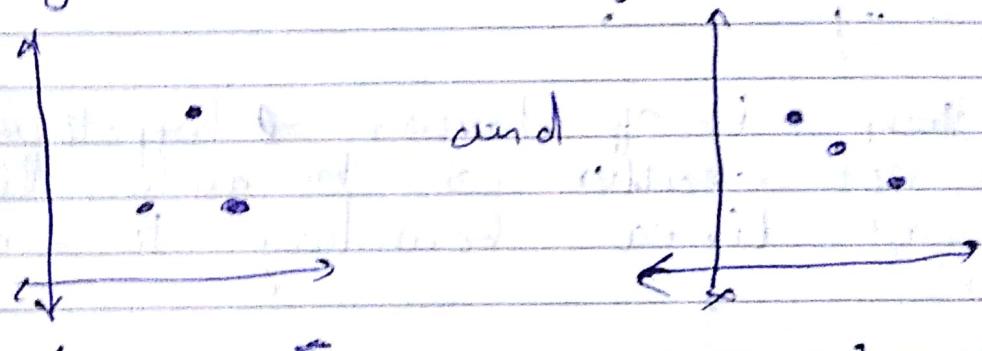
→ For ex.

$H = \{ \text{Linear classifiers in 2D} \}$ , then

$VC(H) = 3$ . → because, we just see that there is no  $H$  which shatters  $S$  with 4 points. in 2D.

6/6/18  
Note

→ Shattering means completely classifying any one set. (Set here means any one arrangements). So ex. of different sets is



corresponding possible classifications.

$$\begin{array}{|c|c|c|} \hline & O & X \\ \hline X & X & X \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & X & O \\ \hline O & X & O \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & X & O \\ \hline X & O & X \\ \hline \end{array},$$

$$\begin{array}{|c|c|c|} \hline & X & O \\ \hline X & X & O \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & O & O \\ \hline O & O & O \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & X & O \\ \hline O & O & O \\ \hline \end{array},$$

$$\begin{array}{|c|c|c|} \hline & O & O \\ \hline X & O & X \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & O & O \\ \hline O & X & O \\ \hline \end{array}.$$

$$\begin{array}{|c|c|c|} \hline & X & O \\ \hline X & X & O \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & X & O \\ \hline O & O & X \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & O & X \\ \hline O & X & X \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline & O & O \\ \hline O & O & O \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & X & X \\ \hline X & X & X \\ \hline \end{array}, \quad \begin{array}{|c|c|c|} \hline & O & O \\ \hline O & O & O \\ \hline \end{array}$$

→ More generally in  $n$ -dimension,

$$VC(\{\text{linear classifier in } n\text{-d}\}) = n+1$$

→ the best result in learning theory :-

+ Theorem - Let  $H$  be given and let  $VC(H) = d$ , then with probability  $1-\delta$ , we've that  $\forall h \in H$ ,

$$|\hat{E}(h) - E(h)| \leq \underbrace{O\left(\sqrt{\frac{d \log m + \log \frac{1}{\delta}}{m}}\right)}$$

Condition  
of uniform  
convergence

→ Thus, with prob.  $1-\delta$ , we also have

$$\hat{E}(h) \leq E(h^*) + O\left(\sqrt{\frac{d \log m + \log \frac{1}{\delta}}{m}}\right)$$

We've already proved this

for 2-D. this is for  
 $n$ -dimensions.

\* Corollary: To guarantee  $E(\hat{h}) \leq E(h^*) + \epsilon$

with probability  $1 - \delta$ , it suffices that

$$m = O_{\epsilon, \delta}(d)$$

→ Assumed

constants while  
writing Big-O notation.

intuition behind this is that, to achieve uniform convergence, we require the number of training examples, is linearly proportional to the VC dimension of the hypothesis class  $H$ .

→ And it turns out for most reasonable hypothesis classes, the VC dimension is very similar to the no. of parameters.

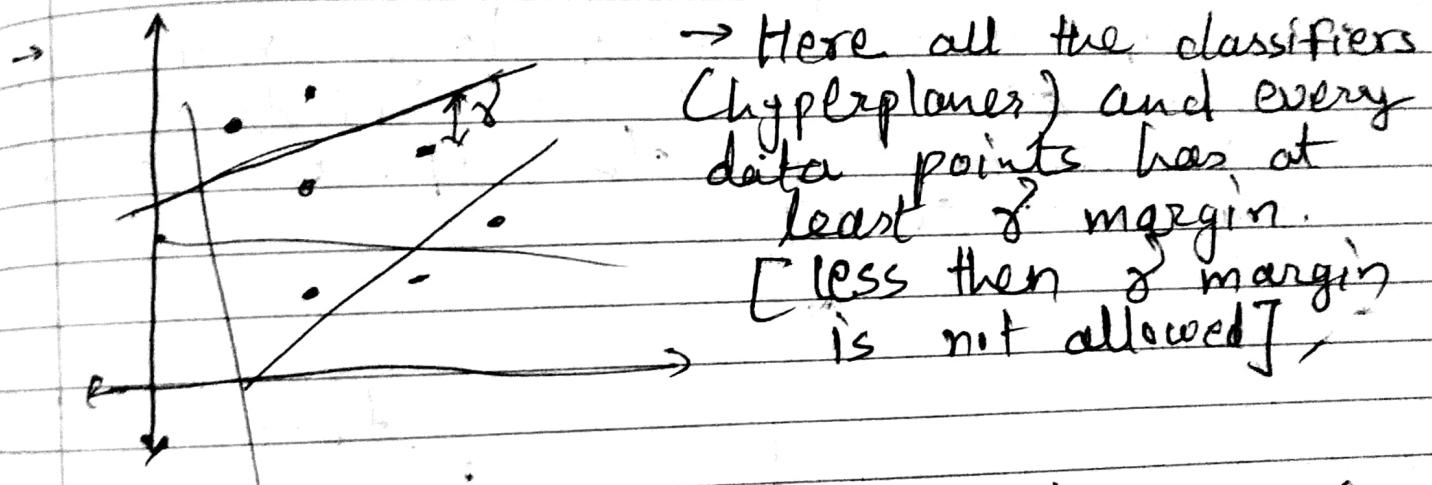
→ For ex. Logistic regression in  $n$ -dimensions has  $n+1$  parameters, and also VC dimension is  $n+1$  [As mentioned previously for linear classifiers].

→ So, no. of training set ex. required is roughly linear to the no. of parameters of model.

[Not true for few strange examples]

## \* How SVM don't overfit?

- In SVMs, using kernels, we map the features to infinite dimension feature space, so it seems like  $VC(H) = \infty$ .
- But it turns out that the class of linear separators with large margin has low VC dimension.



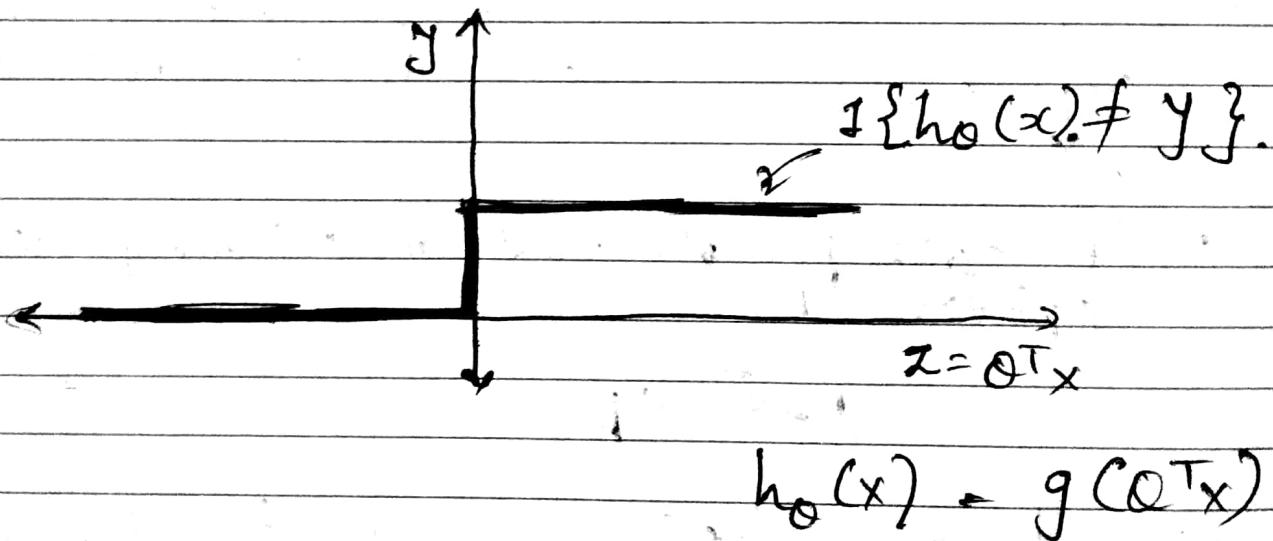
→ If data points lie in within some sphere  $R$ ,

$$\|x^{(i)}\| \leq R$$

→ and considering only the linear separators, which separate the data with margin of at least  $\gamma$ , then

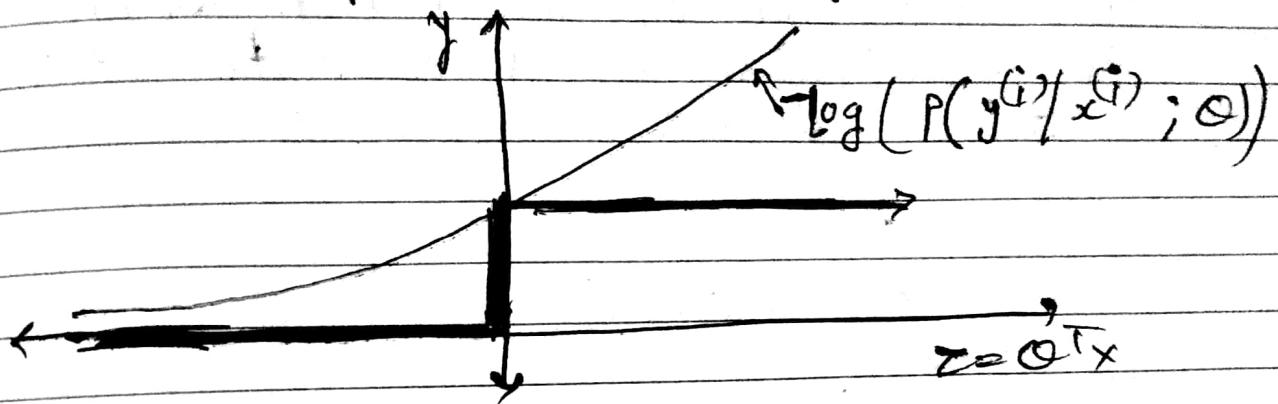
$$VC(H) \leq \left\lceil \frac{R^2}{4\gamma^2} \right\rceil + 1$$

- This is a bound of VC dimensions, which has no dependence on dimension of  $H$ .
- So, even if the data points are in  $\alpha$  dimensions, if we have some class  $H$ , for which margin is high, then VC dimension can be low.
- So, in trying to find large margin, SVM automatically finds lower VC dimensions. class  $H$ . and so it does not overfit.



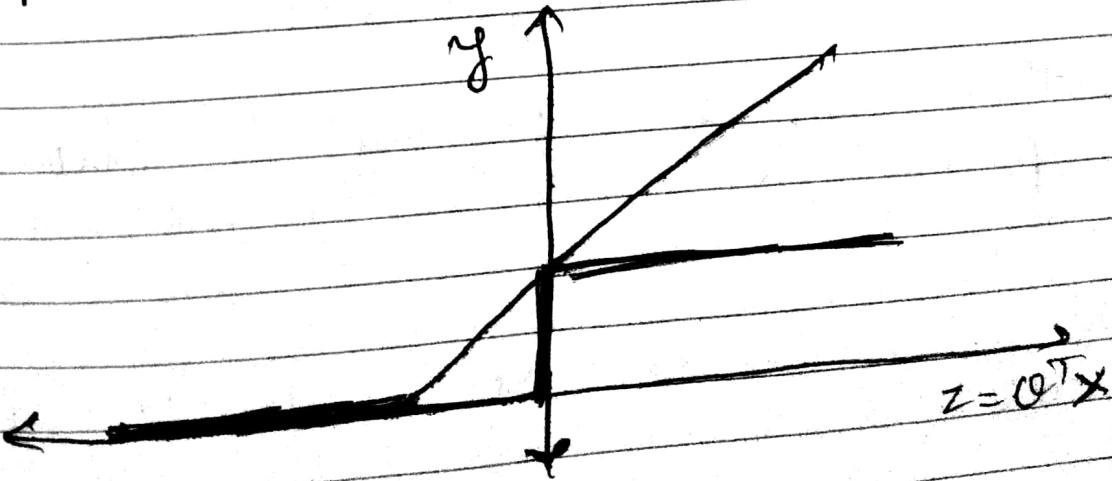
- So, here we want to choose params.  $\theta$  to minimize the above step function [we want  $h_\theta(x) \approx y$ , which is on -ve  $x$ -axis, and that's why we want to minimize it].

- But this step function is non-convex function, so non-convex optimization.
- Both Logistic Regression and SVMs can be viewed as using a rough ~~convex~~ approximation of this problem.



Logistic Regression trying to maximize likelihood [or log likelihood], and so minimize -log likelihood.

Also SVM tries to minimize this step function but with different approx. than logistic regression.



## \* Model Selection:

→ It provides class of methods, to automatically make the trade off between bias and variance.

→ Let's say we are deciding what degree of polynomial to choose to best fit any model.

$$\rightarrow Q_0 + Q_1 x,$$

$$\rightarrow Q_0 + Q_1 x + Q_2 x^2$$

$$\dots$$

$$\rightarrow Q_0 + Q_1 x^1 + \dots + Q_{10} x^{10}$$

→ Another ex. of model selection is to choose

$T$  = bandwidth parameter in locally weighted linear regression.

→ Another ex. is trying to choose param  $C$ , in soft margin SVM.

→ Let's say we've some finite set of models,

$$M = \{ M_1, M_2, M_3, \dots \}$$

may be                       $\theta_0 + \theta_1 x + \theta_2 x^2$   
 $\theta_0 + \theta_1 x$  (linear)      (Quadratic)  
 classifier                  classifier

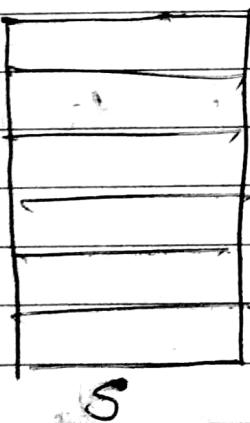
~~One feasible idea is to train the model using each~~

- One terrible idea is to train ~~different~~ all the models above and see which model gives min. error.
  - One method to select the perfect model is Hold-out cross validation.
  - In this we take a training set (all data) and randomly split training set ( $S$ ) into two subsets
    - $S_{train}$  (70%) and
    - $S_{cross-validation}$  (30%)
  - Now, train each model on only  $S_{train}$  and test it on  $S_{cv}$ . and pick model with lowest error of  $S_{cv}$ .

→ This method works well, but in some applications, there are very small amount of training sets or data available, so holding the 30% of total data just for model selection is more costly in such areas.

### \* K-fold cross Validation

→ In this method, take the whole data set  $S$ , divide it in  $K$  pieces and then train all the models on  $K-1$  pieces and then test over remaining 1 piece. And then take average ~~over~~ over the  $K$  results.



→ So for each model, do the above process and take average. Select the model with lowest average.

→  $K=10$  is common choice.

→ Advantage is we only hold  $\frac{1}{K}^{th}$  of data instead of 30% of data for testing.

- But disadvantage is it requires more computation power.  $\Rightarrow$  Bcz each model needed to be trained  $k$  times.
- If we take  $k=m$ , where  $m = \text{no. of examples (total data set)}$ , then it is called leave one out cross validation.
- So it works like take out the 1<sup>st</sup> example and train on rest of the examples and then test on 1<sup>st</sup> one and so on... [very expensive in terms of ~~computation~~ computation].
- So only used when we've too few training examples, like only 15 training examples for ex.

9/6/18

## \* Feature Selection:-

- In many ML problem, features are very high or very high dimensional feature space.
- For ex. in text classification (spam/not-spam), ~~feature~~ no. of features can easily be as large as size of dictionary [50,000 or more].

- So, with such high features, there may be a risk of overfitting.
- If we reduce the no. of features, may be we can reduce the variants of learning algo. and in turn risk of overfitting can be reduced.
- So, ~~the~~ in text classification, there are many such features, which are not relevant to the problem. for ex. if there are words like "and", "is", "not" etc don't tell anything about spam or ~~or~~ not spam. But there, ~~are~~ is also a small set of words which are relevant.
- So, in feature selection, we only select a ~~set of~~ subset of features, which is relevant to specific problem.
- With  $n$  features, there are  $2^n$  possible subsets. [include the feature or not].
- In feature selection, we use different heuristics to search in the ~~the~~ possible feature space of  $2^n$  features.

\* One simple algo. for this is forward search or selection,

, start with  $F = \emptyset$

Repeat {

1) For  $i = 1 \dots n$

try adding feature  $i$  to  $F$ , and evaluate using cross-validation.

2) Set  $F = F \cup \{\text{best feature found in step } 1\}$

}

→ Output the best hypotheses found.

### "Wrapper" Feature Selection

→ This is also called wrapper feature selection. Because it is wrapper over learning algo. bec it calls learning algo. from it.

### Backward search :-

→ Start with  $F = \{1, 2, \dots, n\}$  (all the features) and delete one feature at a time.

- Both this algo. does not guarantee of finding best feature. It is NP-hard problem. But both algo. works fine in practice.
- Forward selection is computationally very intense.

### \* "Filter" Method:-

- For each feature  $i$ , compute some measure of how informative  $x_i$  is about  $y$ .

E.g.  $\text{corr}(x_i, y)$

$$MI(x_i, y) = \sum_{x \in \{0, 1\}} \sum_{y \in \{0, 1\}} P(x_i, y) \log \frac{P(x_i, y)}{P(x_i)P(y)}$$

Mutual Information

Estimate from training Data

$$= KL(P(x, y) || P(x)P(y))$$

KL divergence :- Shows how diff. two probability distribution are.

Kullback-Leibler divergence

→ So, if KL divergence is zero,

$$\text{then } P(x, y) = P(x) P(y)$$

that means  $x$  and  $y$  are independent  
and they are less mutual informative.

→  $x$  and  $y$  are ~~non-independent~~ non-independent  
then it's assumed that  $x$  can tell  
something about  $y$  or vice-versa.

so, they are highly mutual informative.

→ Having chosen above ~~measure~~ measure,  
pick top  $k$  features, with largest  
 $\text{Corr}(x_i, y)$  or  $\text{MI}(x_i, y)$ .

→ And for choosing  $k$ , we can use Cross-validation.

## Lecture - II

### \* Bayesian statistics and Regularization

- In previous lecture learning theory, we saw that those theories tries to simplify the model by selecting simple models.
- For ex, in feature selection, it was about to reduce the no. of features and hence simplify the model.
- Regularization is another way to prevent overfitting and it also keep model complexity same as previously.
- In linear Regression,

maximum likelihood

$$\max_{\Theta} \prod_i P(y^{(i)} | x^{(i)}; \Theta)$$

- This is called frequentist procedure.
- The alternative of this method of estimating is Bayesian statistics.
- $P(\Theta)$  - prior

e.g. prior on  $\Theta$  can be gaussian distribution as

$$\boldsymbol{\theta} \sim N(\boldsymbol{0}, \tau^2 \mathbf{I})$$

→ let data set is,

$$S = \{(x^{(i)}, y^{(i)})\}_{i=1}^m$$

→ then in Bayesian statistics, we'll calculate posterior probability of  $\boldsymbol{\theta}$  given  $S$ .  $P(\boldsymbol{\theta}|S)$

$$\rightarrow P(\boldsymbol{\theta}|S) \propto \left( \prod_{i=1}^m P(y^{(i)}|x^{(i)}; \boldsymbol{\theta}) \right) P(\boldsymbol{\theta})$$

$T_{\text{posterior}}$

→ To make a new prediction on new input  $x$ ,

$$P(y|x, S) = \int_{\boldsymbol{\theta}} P(y|x, \boldsymbol{\theta}) P(\boldsymbol{\theta}|S) d\boldsymbol{\theta}$$

$$E[y|x, S] = \int y P(y|x, S) dy$$

here  $\boldsymbol{\theta}$  is considered as random variable, [see it's now separated by comma instead of ;].

MAP = Maximum A posteriori estimation.

Page No.		
Date		

- The above integral is commonly hard to calculate when  $\theta$  is very high dimensional.
- so, instead of computing above, we will do following.

$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} P(\theta | s)$$

$$= \operatorname{argmax}_{\theta} \left( \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta) \right) P(\theta).$$

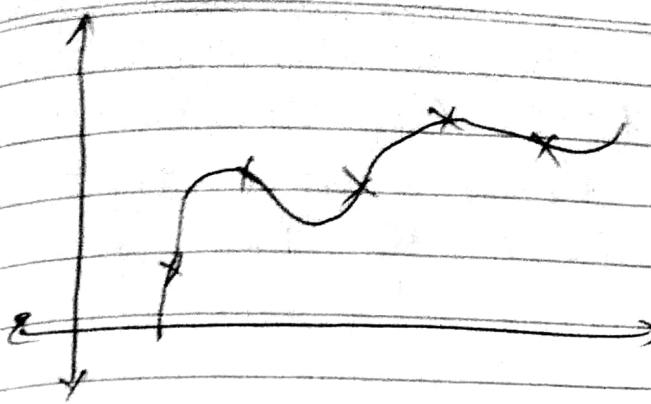
- And to make prediction,

$$\text{product } h_{\hat{\theta}_{MAP}}(x) = \hat{\theta}_{MAP}^T x$$

- One intuition is that if our  $\theta$  is gaussian distributed,

$\theta \sim N(0, \Sigma^{-1})$ , then it's probability mass is concentrated near zero, so it is same as setting feature values to zero and eliminate them, as done by feature selection previously.

- So, in max. likelihood procedure, if we fit higher order polynomials to small data sets, it will overfit.



- But in Bayesian regularization, the above curve becomes smoother as we decrease  $T$ . So even we select higher order polynomial [Complex model] it still doesn't overfit.
- Another intuition is that for ex. in linear regression, max. likelihood minimize the following cost function.

$$\min_{\theta} \sum_i \|y^{(i)} - \theta^T x^{(i)}\|^2$$

- But this Bayesian minimizes
- $$\min_{\theta} \sum_i \|y^{(i)} - \theta^T x^{(i)}\|^2 + \lambda \|\theta\|^2$$

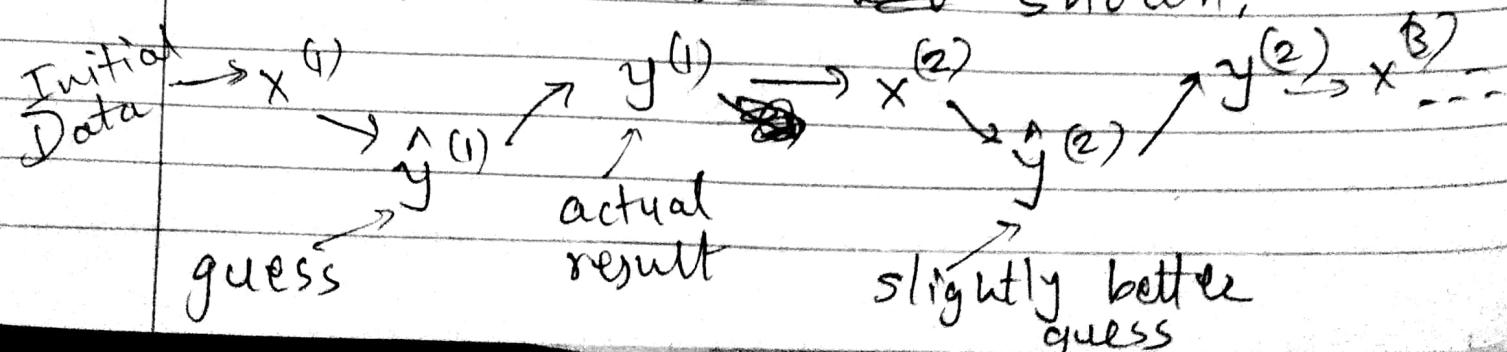
→ So, we add penalty for being  $\theta$  large, so it tries to  $\theta$  predict as small as possible and also minimize the cost.

→ This type of regularization is used generally in text classification, where we have very high feature vector and low training examples.

## \* Online Learning:

### \* Online Learning:

- So far we've seen batch learning algorithms, where we've given a training set to train the model and another set of data known as test set for testing or evaluating the model.
- But in online learning, we make predictions, even while we are learning (training the model).
- Let's take example of classification problem.
- Initially, without seeing any data we'll first guess the result [Just blind guess], then the original result will be ~~shown~~ shown,



→ Here we care about  
Total online error:-

$\sum_{i=1}^m I\{\hat{y}^{(i)} \neq y^{(i)}\}$ , which is total no.  
of wrong guess  
made.

→ One good example, where Online learning algo. gives good result is perceptron learning algo.

Initialize  $\Theta = 0$

After  $i^{th}$  Example, update

$$\Theta = \Theta + \alpha (y^{(i)} - h_\Theta(x^{(i)})) x^{(i)}$$

## Lecture-12

Page No.	166
Date	16/6/18

### \* Unsupervised Learning:-

- Given a dataset without any labels and from this, it's job of learning algo. to find out patterns or structure in the given data set.
- For example clustering ~~is~~ algorithm which is used in clustering group of customers for business, in biology for clustering genes.
- "news.google.com" uses such clustering algo. to ~~as~~ categorizes news in different ~~category~~ categories.
- Image segmentation is also an ex. of clustering algo.

### \* K-means clustering Algo:-

Input will be unlabeled data set

$$\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\} \in \mathbb{R}^n$$

Note that there is no class labels  $y$  in the datasets because this is input to unsupervised learning.

1<sup>st</sup> step:- Initialize cluster centroids.

$$u_1, u_2, \dots, u_k \in \mathbb{R}^n$$

2) Then repeat until convergence:-

i) set  $c^{(i)} = \arg \min_j \|x^{(i)} - u_j\|$

ii)  $u_j = \frac{\sum_{i=1}^m \mathbf{1}_{\{c^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m \mathbf{1}_{\{c^{(i)}=j\}}}$

→ Refer the animation shown in video for perfect understanding.

\* Does this algo. converge?

→ Answer is Yes.

→ If we define Distortion function

$$J(c, u) = \sum_{i=1}^m \|x^{(i)} - u_{c(i)}\|^2$$

↓  
 data point      ↑  
 i                  centroid corresponding  
 to data point  
 i

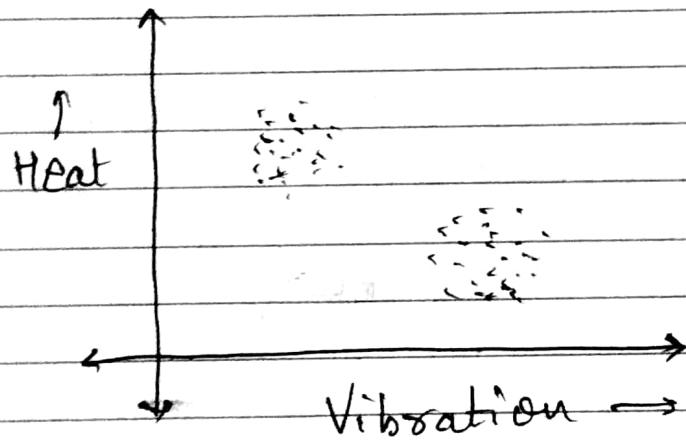
→ Then K-means is coordinate ascent on  $J(c, u)$ .

## \* How to choose no. of clusters?

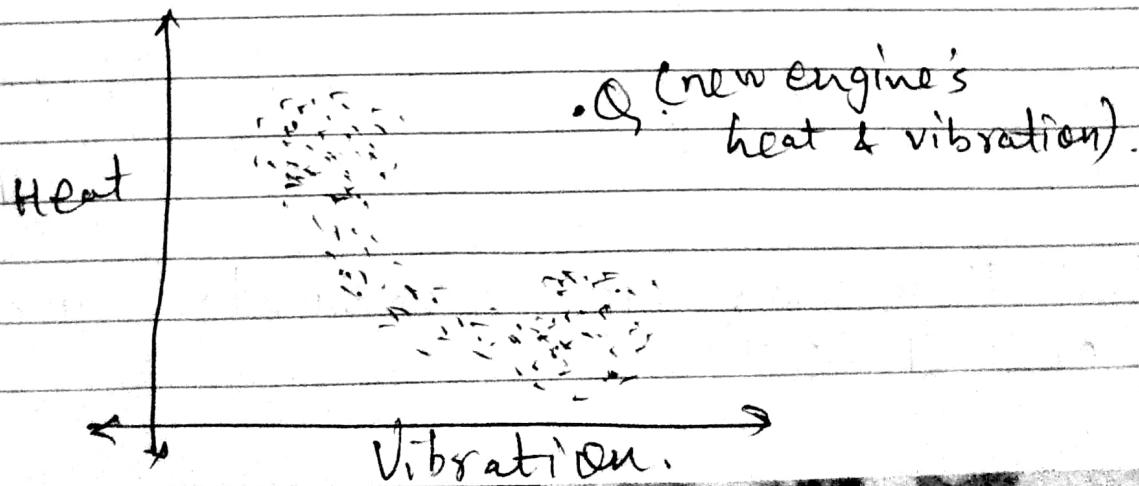
→ It's try-error approach. Pick any suitable number and then try the algorithm and pick the no. of clusters for which algo. gives best result.

## \* Density Estimation:

→ Let's say, in an air-craft engine producing company, we've following data for different engines.



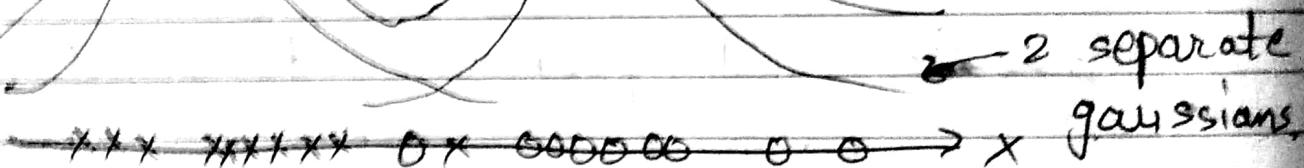
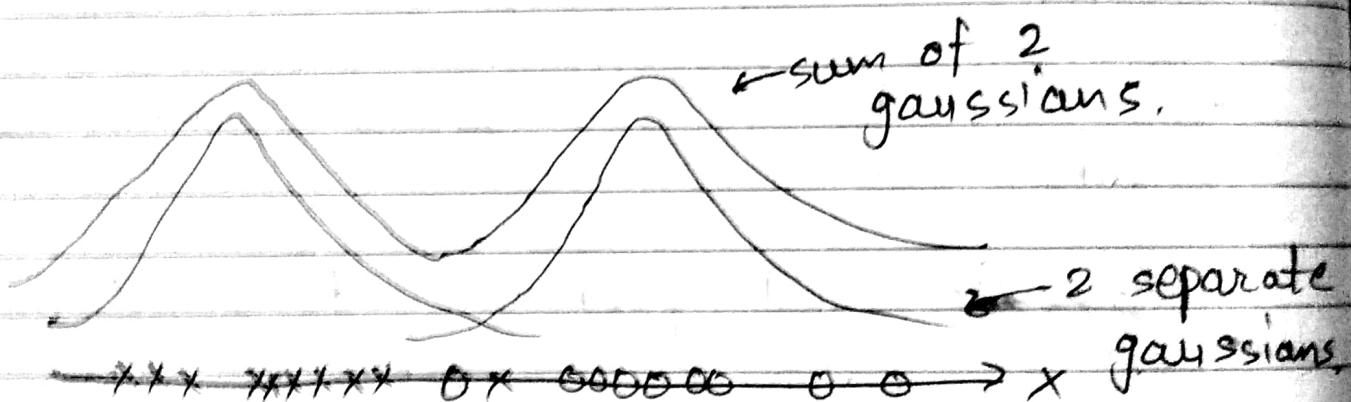
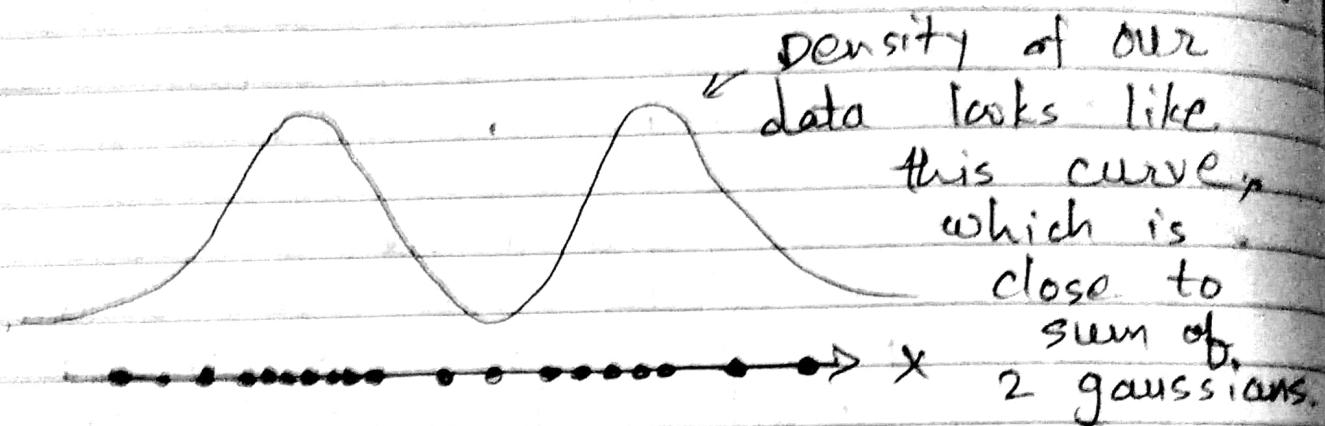
→ Then we can apply clustering algo. to find 2 clusters in above data set. But if we've following data,



- And from this we would like to find joint distribution of amount of heat produced corresponding to amount of vibration.
- Because for any new engine, we would like to decide that it needs further ~~reco~~ inspection or not.
- So, can we build a  $f^n$   $P(x)$  from the above data from which we can say for every new engine in assembly line, whether it needs inspection or not.
- This problem is called ~~error~~ anomaly detection.
- And for such problem, take dataset ~~data~~  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$  and build the model
- If  $P(x)$  is very low, then we can say that it is an ~~normal~~ anomalous example.
- $P(x)$  ~~used~~ here is any distribution function, but it is not like any of the standard distribution functions like gaussian ~~etc~~.

→ So, we want to come up with a model, which can be used to measure such densities.

Ex In 1 D, let's our dataset looks like following.



→ But the problem with our data set is we don't know the actual labels of data.

→ So, now we want to come up with an algo. ~~to~~ to fix the mixture of such gaussians even when we don't have labeled data.

→ There's a latent (hidden) unobserved random variable  $z$ . And  $x^{(i)}, z^{(i)}$  have a joint distribution:-

$$P(x^{(i)}, z^{(i)}) = P(x^{(i)} | z^{(i)}) p(z^{(i)})$$

$z \sim \text{Multinomial}(\phi)$  [Bernoulli for  $z$  Gaussians]  
 $(\phi_j \geq 0, \sum_j \phi_j = 1)$ .

so,

$$P(x^{(i)} | z^{(i)}) = j \sim N(\mu_j, \Sigma_j)$$

→ Same as Gaussian Discriminant analysis, the only diff. is that the labels variable  $y$  is now replaced with latent  $z$ .

→ Now, to link this to GDA, let's assume we know  $z_i$ 's.

→ Then we could write max.<sup>log</sup> likelihood as,

$$l(\phi, \mu, \Sigma) = \sum_{i=1}^m \log p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)$$

$$\phi_j = \frac{1}{m} \sum_{i=1}^m I\{z^{(i)} = j\}$$

$$\mu_j = \frac{\sum_{i=1}^m I\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m I\{z^{(i)} = j\}}, \text{ etc.}$$

## Maximization

\* EM (Expectation ~~and~~ Maximization) Algo.

Repeat {

E-step: Guess values of  $z_i$ 's

$$\text{Let } w_j^{(i)} = P(z^{(i)}=j | x^{(i)}, \phi, \mu, \Sigma)$$

This is probability that data point  $x^{(i)}$  come from Gaussian number  $j$ . ( $z_j$ )

$$= \frac{P(x^{(i)} | z^{(i)}=j) P(z^{(i)}=j)}{\sum_{l=1}^k P(x^{(i)} | z^{(i)}=l) P(z^{(i)}=l)}$$

from bayes theorem

$$\propto \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp\left((x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \cdot \phi_j$$

$$w_j = \frac{1}{\sum_{l=1}^k \frac{1}{(2\pi)^{d/2} |\Sigma_l|^{1/2}}} \exp\left((x^{(i)} - \mu_l)^T \Sigma_l^{-1} (x^{(i)} - \mu_l)\right) \cdot \phi_l$$

M step :- (Maximization step) :-

$$\phi_j = \frac{1}{m} \sum_{i=1}^m w_j^{(i)}$$

$$\mu_j = \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}$$

$$\varepsilon_j = \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j) (x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}}$$

?

→ Note that the formulas of  $\phi_j$  and  $w_j$  are similar to GDA. Also in GDA we were using ~~only~~ ~~one~~ same ~~cov~~ covariance matrix for all gaussians, but here it is different for ~~for~~ each gaussian.

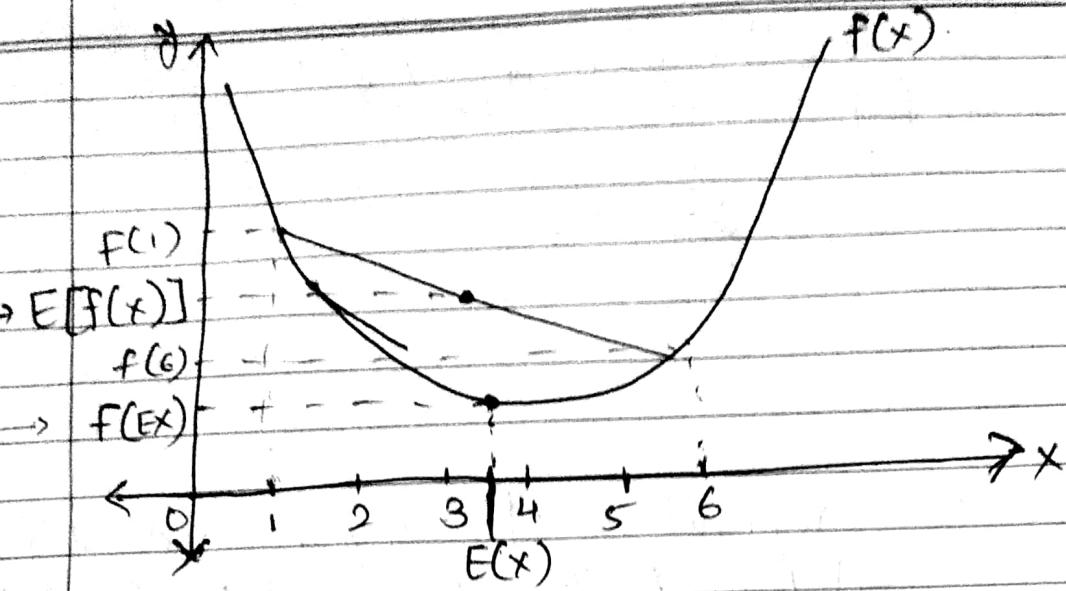
→ This was special case of more general EM algorithm. For deriving General form of EM, following knowledge is required.

\* Jensen's inequality:

→ Let  $f$  be a convex function, (e.g.  $f''(x) \geq 0$ ) and let  $x$  be a random variable. Then,

$$f(E[X]) \leq E[f(X)]$$

↑  
expectation  
of  $X$ .      ↑  
expectation  
of  $f(X)$ .



Expected value of  $x = \text{avg of all } x$   
 $= \frac{1+2+3+4+5+6}{6}$   
 $= 3.5$

→ and same can be done for  $E[f(x)]$   
 and the result can be proved.

→ further, if  $f''(x) > 0$  (Strictly convex),  
 then,  
 $E[f(x)] = f(E[x]) \Leftrightarrow \underbrace{x = E(x)}$  with prob. 1.  
 this means when  $x$  is constant.

→ If  $f'' \leq 0$  [non convex function], then,

$$f(E[x]) \geq E[f(x)]$$

## Generalized EM:

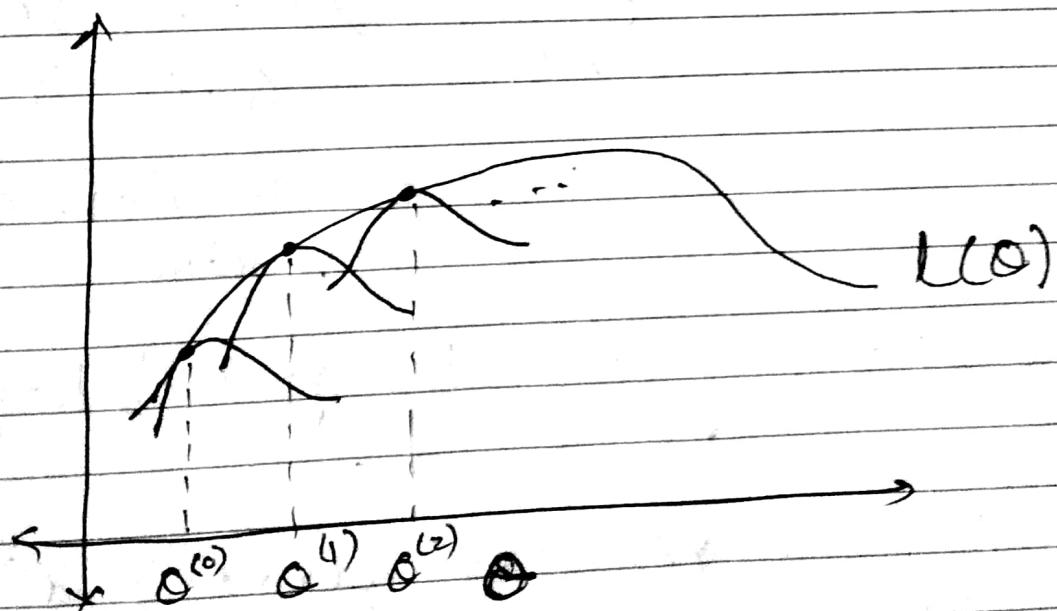
- We've some model for  $P(x, z; \theta)$ , we observe only  $x$ .
- Our goal is to maximize

$$l(\theta) = \sum_{i=1}^m \log P(x^{(i)}; \theta)$$

log likelyhood  
of param.  
of model.

$$= \sum_{i=1}^m \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)$$

- What this algo. does is, initially starts with  $\theta^{(0)}$  and constructs a lower bound to  $l(\theta)$  and then maximizes it.



- We want to maximize  $l(\theta)$ , so,

$$\max_{\theta} l(\theta) = \max_{\theta} \sum_i \log P(x^{(i)}; \theta)$$

$$= \sum_i \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)$$

$$= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

$Q_i$  is probability distribution function, such that,

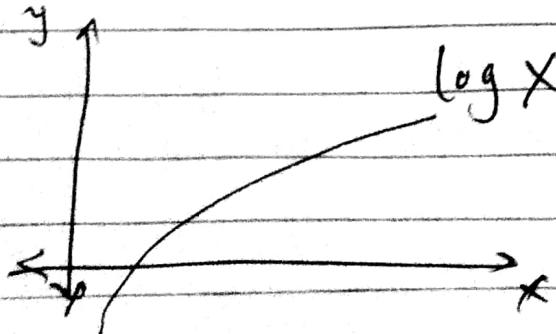
$$Q_i(z^{(i)}) \geq 0, \text{ and } \sum_{z^{(i)}} Q_i(z^{(i)}) = 1$$

$$= \sum_i \log E_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

[Refer to def. of expectation]

[ $Q_i(z^{(i)})$  is also the probability of  $z^{(i)}$ ].

$\Rightarrow \log x$  is concave [non-convex].  $f'' < 0$ .



→ so, from Jensen's inequality,

$$\log E[x] \geq E[\log x].$$

∴ so,

$$\sum_i \log \mathbb{E}_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \geq$$

$$\sum_i \mathbb{E} \left[ \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

$$= \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

→  $z \sim P$ , and a function  $g(z)$ , the

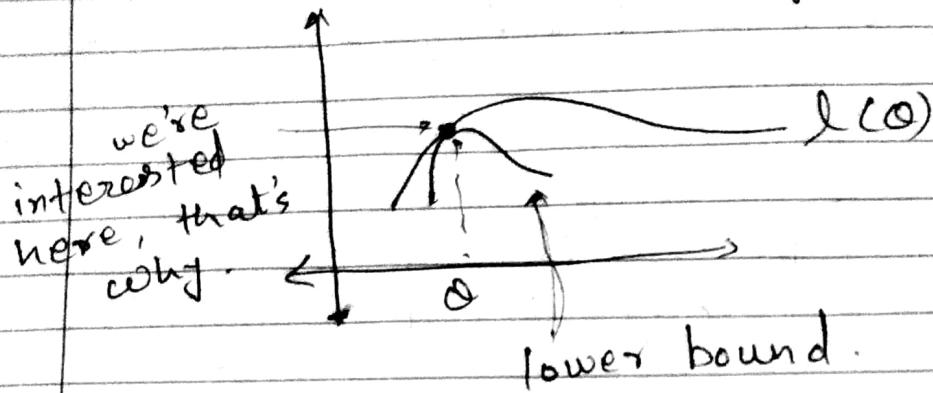
$$E[g(z)] = \sum_i P(z) g(z)$$

→ we have derived the lower bound of  $l(\theta)$ , so,

$$l(\theta) \geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

→ now, to make above inequality to be equality, we need to choose  ~~$Q_i$~~  such distribution  $Q_i$ .

- Why we need equality?  
 → Refer to the following intuition.



- So, from, Jensen's inequality, for equality cond<sup>n</sup>,

$$\frac{P(x^{(i)}, z^{(i)}; \theta)}{\sum Q_i(z^{(i)})} = \text{constant} \quad (\text{some value for all values of } z^{(i)})$$

$$\text{so, } Q_i(z^{(i)}) \propto P(x^{(i)}, z^{(i)}; \theta)$$

and  $Q_i$  is probability distribution,  
 so,

$$\sum_{z^{(i)}} Q_i(z^{(i)}) = 1$$

$$\rightarrow Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)}; \theta)}{\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)}$$

$$= \frac{P(x^{(i)}, z^{(i)}; \theta)}{P(x^{(i)}; \theta)}$$

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

→ So, EM algo given by :-

Repeat {

E-step :-

$$\text{choose/set } Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

M-step :-

$$\theta = \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

}.