

24/5/18

# Machine Learning | Coursera | Stanford

## \* WEEK 1 :

Click stream data : web click data being collected for ML purpose.

Note SVM allows the computer to deal with infinite number of features.

e.g. of unsupervised learning - (cocktail party problem algorithm)

Microphone recordings with two people speaking simultaneously. The unsupervised algo will distinguish the two voices as the diff microphone recordings have diff amplitudes of the two speaking people thus giving two separate outputs with individual speaking in recording 1 as output 1 and output 2 with individual 2's voice separated out.

Training Set

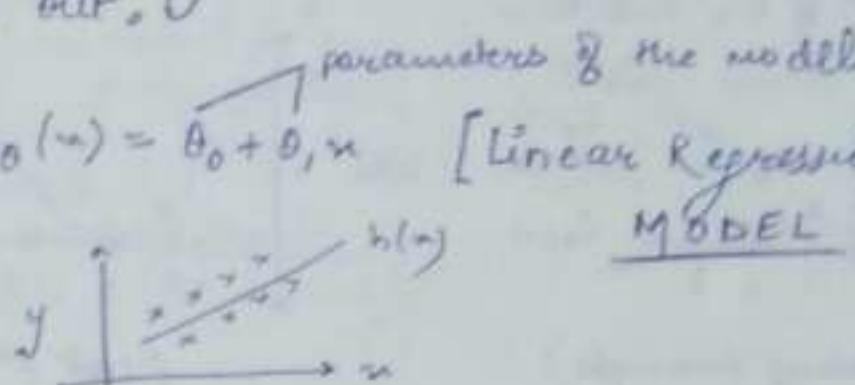
Learning Algo

I/P  $\rightarrow$   $\boxed{w}$   $\rightarrow$  O/P  
(hypothesis)

a fn to map i/p to corresponding o/p

Hypothesis is the standard terminology used by people in ML }

$$\text{e.g. } h_0(x) = \theta_0 + \theta_1 x \quad [\text{Linear Regression}]$$

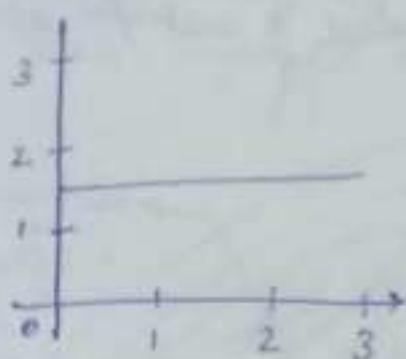


fitting linear fn on the training dataset

Linear regression is one variable

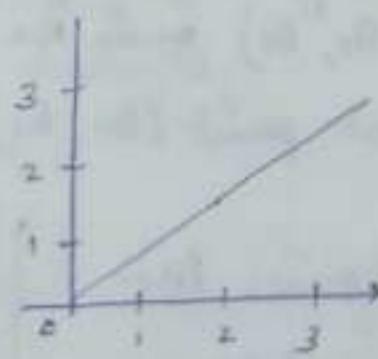
① = Univariate Linear Regression

diff  $\theta_0$  and  $\theta_1$  lead to different hypothesis fn,  $h_0(x)$ .



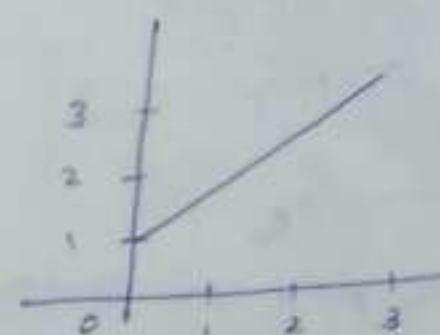
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

$$\text{for } h_0(x) \\ = \theta_0 + \theta_1 x,$$

We choose  $\theta_0, \theta_1$  so that  $h_0(x)$  is close to  $y$  for our training e.g. ( $x, y$ )

$$\text{minimize } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

i.e. finding  $\theta_0$  and  $\theta_1$  to minimize this expression.

where  $m = \text{no. of training examples}$

overall objective fn for linear regression.

i.e. minimize  $J(\theta_0, \theta_1)$

$\theta_0, \theta_1$

where  $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$  is our goal to minimize.

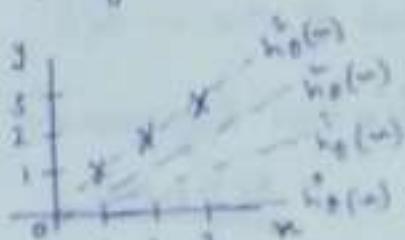
cost fn for linear regression

(Squared error fn)

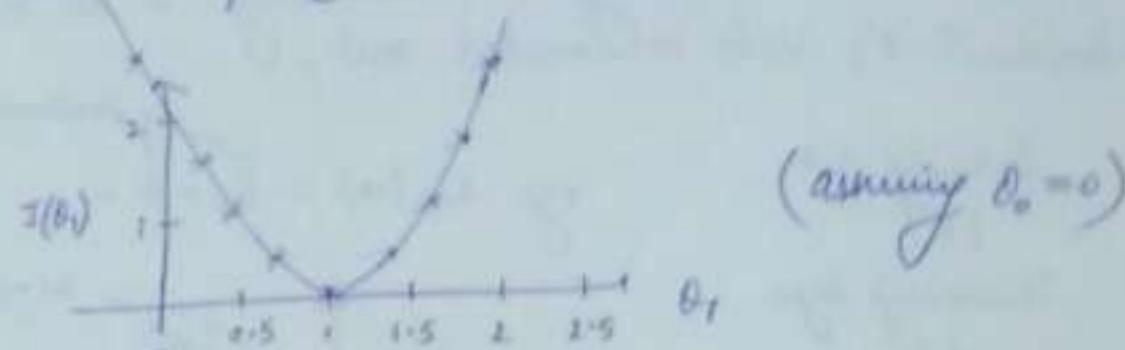
$h_\theta(x)$ : for fixed  $\theta_0$  and  $\theta_1$ , this is a fn of  $x$ .

$J(\theta)$  is fn of  $\theta_0$  and  $\theta_1$ .

e.g.



(training example)

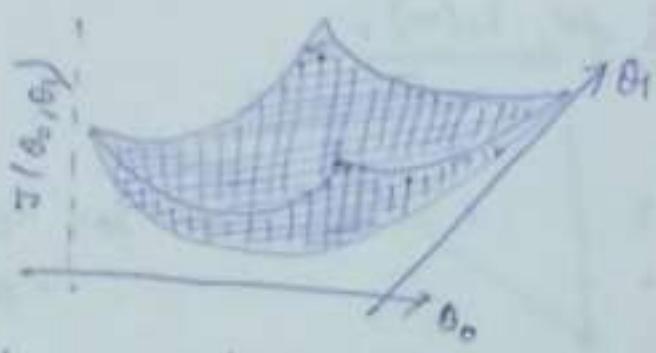


(cost fn values)

Each value of  $\theta_1$  corresponds to a different hypothesis  $h_\theta(x)$  or to a different straight line fit on the left and for each value of  $\theta_1$ , different value of  $J(\theta_1)$  can be derived.

Above  $J(\theta)$  graph  $\Rightarrow \theta_1 = 1$  is the best parameter value giving min  $J(\theta_1)$

when both  $\theta_0$  and  $\theta_1$  are considered,  $J(\theta_0, \theta_1)$  is a 3D plot like:

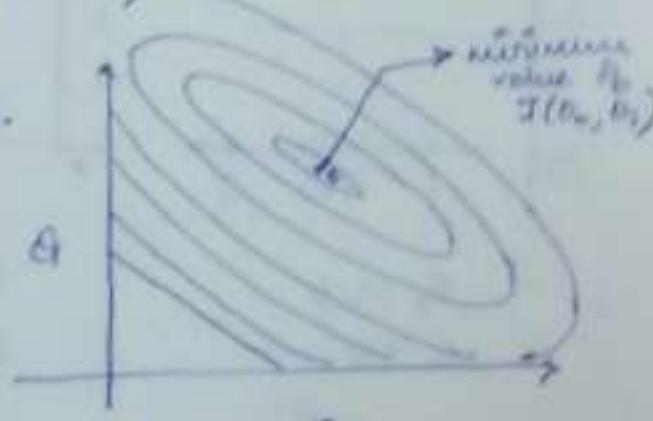


The height of the graph at certain  $(\theta_0, \theta_1)$  gives the value of cost fn at that point  $(\theta_0, \theta_1)$

We will use contour plots or contour figures like.

The ovals show the set of points that takes on the same value for  $J(\theta_0, \theta_1)$ . Points lying on the same oval have the same value.

gives  $(\theta_0, \theta_1)$  for global minimum.



So, we need an algo to automatically find value of  $\theta_0$  and  $\theta_1$  that minimizes cost function  $J$ . It is called gradient descent. It is used to minimize other  $fN$  as well not just cost  $fN$  for linear regression. Also called batch gradient descent.

25/5/18

In this algo, we start at some  $\theta_0, \theta_1$  (say 0) and keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we get minimum or local minimum. Gradient descent produces different minimum based on different starting points.

repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0 \text{ and } j=1)$$

Step 1: temp 0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

assignment operator (assigns RHS to LHS)

Step 2: temp 1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

while  $a = b$  is truth assertion

Step 3:  $\theta_0 := \text{temp 0}$

(simply asserts that  $a = b$ )

Step 4:  $\theta_1 := \text{temp 1}$

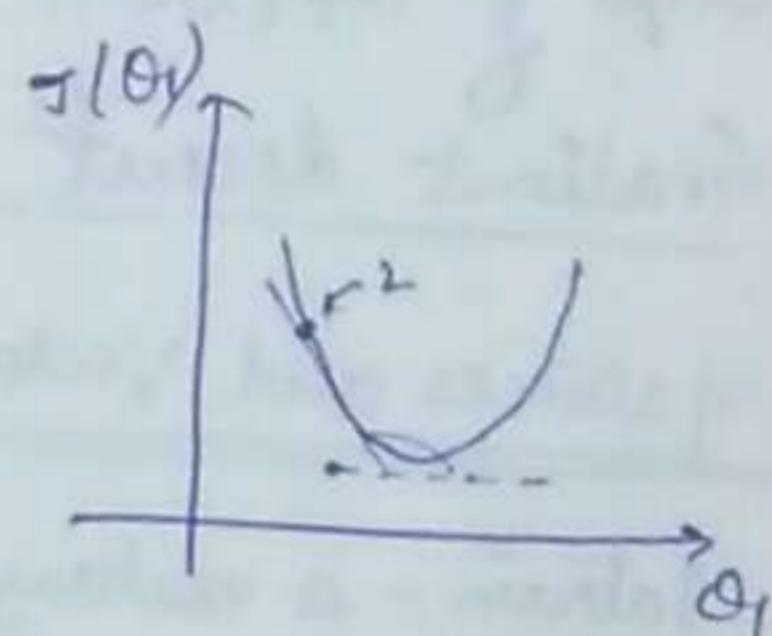
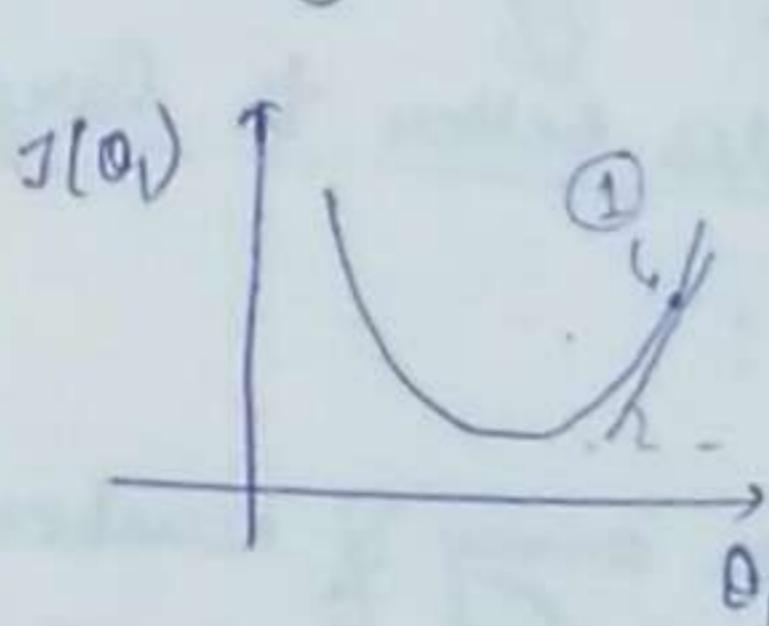
Hence, we can write  $a := a + 1$   
but not  $a = a + 1$ .

can't do step 3 before

Step 2 bcoz otherwise in

$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$  updated  $\theta_0$

learning rate  
(decides how big step to  
be taken while descending)



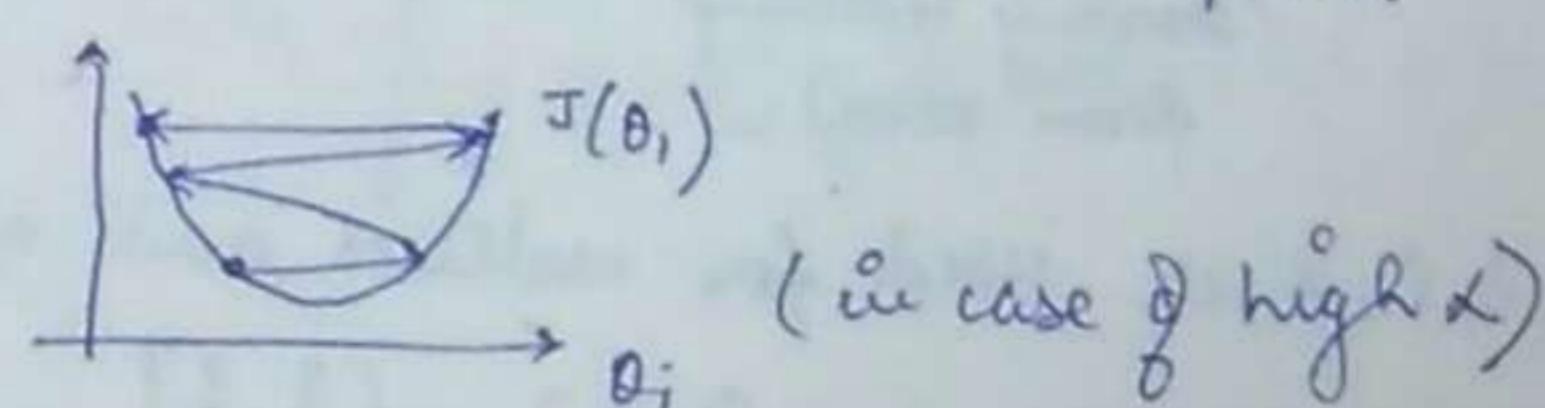
Note If  $\alpha$  is too small, gradient descent can be slow, if its too large it can overshoot the minimum hence failing to even converge or diverge.

e.g. 1

In this case  $\theta_1$   
will  $\downarrow$  bcoz of 1 being  
the starting point.

e.g. 2

In this case  $\theta_1$   
will  $\uparrow$  bcoz of 2  
being the starting  
point.



If  $\theta_1$  is initially at local minimum, gradient descent will leave  $\theta$  changed to nothing different but original  $\theta$ , only as  $\frac{d}{d\theta_1} J(\theta_1) = 0$ .

Gradient descent can converge to a local minima, even  $\alpha$  fixed bcoz gradient " will automatically take smaller steps and so no need to  $\downarrow \alpha$  over time.

Applying gradient descent to minimize squared error cost fn.

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} \quad \text{as } h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)} \rightarrow \text{slope of f}$$

$$\text{and we know: } \theta_j = \theta_j - \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \quad \{ j=0, 1 \}.$$

Convex fn: a bowl shaped fn and so has just one global optimum instead of having any local optima.

The above algorithm is also called "Batch Gradient Descent". It refers to the fact that we're looking at all of the training examples. So, in gradient descent when computing the derivatives we are computing the sums over our 'm' training examples and so term refers to the fact that we are looking at the entire batch of training examples.

There are methods that solves for  $\min J(\theta_0, \theta_1)$  without needing multiple steps of iterations like in gradient descent.

Gradient descent scales better to larger data sets than normal eq method

Matrices and Vectors:

Matrix - a rectangular array of numbers.

e.g.  $\begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$  is  $\mathbb{R}^{3 \times 2}$

Dimension of matrix - (no. of rows)  $\times$  (no. of columns)

Vector - An  $(n \times 1)$  matrix i.e. matrix with just 1 column. e.g.  $\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$  is a 4-D vector  
↑ index ordered vector  
(starting indexing from zero)

A, B, C, X, etc used for matrices and a, b, c, n for vectors representation commonly.

Scalar multiplication:  $3 \times \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix}$  (it is commutative)

e.g.  $h_\theta(x) = -40 + 0.25x$

$$\frac{1}{n} \sum_{i=1}^n \begin{bmatrix} 21 \\ 14 \\ 15 \\ 85 \end{bmatrix}$$

Then  $h_\theta(x)$  can be calculated as:  $\begin{bmatrix} 1 & 21 \\ 1 & 14 \\ 1 & 15 \\ 1 & 85 \end{bmatrix} \times \begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$  for (in one go)

$$\therefore \text{Prediction} = \text{Data Matrix} \times \text{parameters}$$

(more computationally efficient.)

Matrix matrix multiplication is a step to the gradient descent alternative approach for calculating  $\theta_0$  and  $\theta_1$ .

Note  $A \times B \neq B \times A$  (matrix multiplication is not commutative)  
 $(A \times B) \times C = A \times (B \times C)$  (• • • associative)

$I_{mn}$ : Identity matrix of dimension  $m \times n$ . (has 1's along the diagonals)  $\begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}$

$A \times B = B \times A$  when  $B$  is Identity matrix (but of diff. dimension)

Matrix inverse: If  $A$  is an  $m \times m$  matrix, and if it has an inverse,

$$A(A^{-1}) = A^{-1}A = I \quad \left\{ \begin{array}{l} A \text{ is a square matrix and only} \\ \text{square matrices have inverses} \end{array} \right.$$

Just like 0, matrix  $A$  if a zero matrix has no inverse. In ML, models not having inverses can be considered close to zero.

Singular / Degenerate Matrix: Matrices that don't have an inverse.

If  $B = A^T$  then  $B_{ij} = A_{ji}$

29/5/18

### Multivariate Linear Regression

$x^{(i)}$  = input (features) of  $i^{th}$  training example.

$x_j^{(i)}$  = value of feature  $j$  in  $i^{th}$  .. .. ..  
 $n$  = no. of features.

$$\text{Now}, h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad \left\{ \begin{array}{l} x_0 = 1 \\ \text{adding an additional feature} \end{array} \right.$$

$$\text{Now}, x \approx \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \theta^T x \quad \{(n+1) \times 1 \times (n+1)\}$$

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = J(\theta)$$

### Gradient descent

$$\text{Repeat } \left\{ \theta_j = \theta_j - \frac{\partial}{\partial \theta_j} J(\theta) \right\}$$

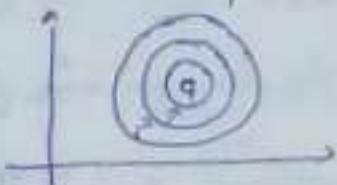
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \left\{ \begin{array}{l} \text{simultaneously update } \theta_j \text{ for} \\ j = 0, \dots, n \end{array} \right\}$$

Note Gradient descent converges more quickly when different features take on similar ranges of values.



In this case where  $\theta_0$  and  $\theta_2$  are in different ranges, gradient descent has a hard time reaching the global minima.

If  $x_1$  and  $x_2$  i.e. features are in same range:



It is easier for gradient descent to reach global minimum.

for this mean normalization can be done i.e

replacing  $x_i$  with  $\frac{x_i - \bar{x}_i}{s_i}$  to make features have approx 0 mean.

e.g. no. of bedrooms

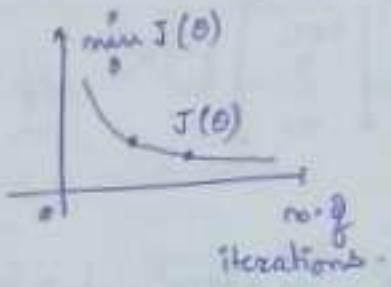
1	$\frac{1-2}{5} = -0.2$
2	$\frac{2-2}{5} = 0$
3	$\frac{3-2}{5} = 0.2$
4	$\frac{4-2}{5} = 0.4$
5	$\frac{5-2}{5} = 0.6$

Normalized values.

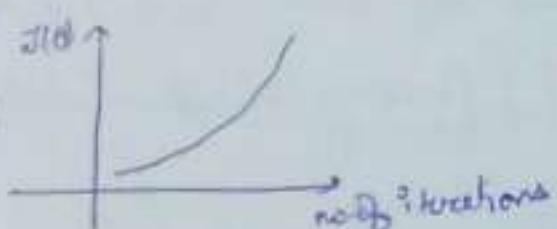
} feature scaling

30/5/18

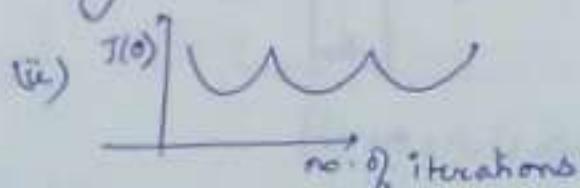
This is called feature scaling to converge gradient descent in much lesser iterations. Gradient descent basically finds  $\theta$  for which  $J(\theta)$  is minimum.



When gradient descent work properly,  $J(\theta)$  goes with every iteration. The plot helps us judge if the gradient descent is converging. Automatic convergence test also helps us do so by checking if  $J(\theta)$  goes by less than  $\epsilon$  in 1 iteration. but determining this  $\epsilon$  for different applica is difficult



- gradient descent working improperly.
- need to +  $\alpha$ .



→ need to use smaller  $\alpha$ .

Notes:  $\alpha$  should not be too small, as then we would need many iterations to get minimum  $\theta$ .

Note If  $\alpha$  is too large, slow convergence is possible but  $J(\theta)$  may not ↓ on every iteration.

So, best way is to plot  $J(\theta)$  for different values of  $\alpha$  against the no. of iterations.  $\alpha$  that leads to convergence of  $J(\theta)$  should be selected.

→ Polynomial regression:

e.g.  $h_\theta(u) = \theta_0 + \theta_1 u + \theta_2 u^2 + \theta_3 u^3$  (cubic fit to my data)

$$\begin{array}{c} \downarrow \\ u_1 \end{array} \quad \begin{array}{c} \downarrow \\ u_2 \end{array} \quad \begin{array}{c} \downarrow \\ u_3 \end{array}$$

$\underbrace{\hspace{1cm}}$  features.

e.g.  $h_\theta(u) = \theta_0 + \theta_1 u + \theta_2 \sqrt{u}$

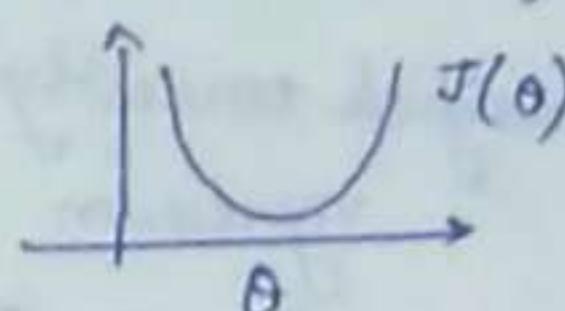
• Normal eqn - an alternate to gradient decent:

method to solve for  $\theta$  analytically, so that  $\theta$  can be obtained in 1 go instead of repetitive iterations as in gradient decent.

e.g. if 1D ( $\theta \in \mathbb{R}$ )

$$J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{d}{d\theta} J(\theta) = 0 \Rightarrow \text{solve for } \theta.$$



In general:  $\frac{\partial}{\partial \theta_j} J(\theta) = 0$  (for every  $j$ )

and solve for  $\theta_0, \theta_1, \dots, \theta_n$ .

$$u^{(i)} = \begin{bmatrix} u_0^{(i)} \\ u_1^{(i)} \\ \vdots \\ u_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

, so,  $X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$  = design matrix.

e.g. if  $u^{(i)} = \begin{bmatrix} 1 \\ u_1^{(i)} \end{bmatrix}$ ,  $X = \begin{bmatrix} 1 & x_1^{(i)} \\ 1 & x_2^{(i)} \\ \vdots & \vdots \\ 1 & x_m^{(i)} \end{bmatrix}_{(m \times 2)}$  and  $y = \begin{bmatrix} y^{(i)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

Feature scaling isn't necessary when we use Normal eqn.

Gradient Decent

- i) Need to chose  $\alpha$
- ii) Many iterations needed
- iii) works well even when 'n' (features) is large.
- iv) when  $n = 10^4$ , we switch to gradient descent.

Normal eqn

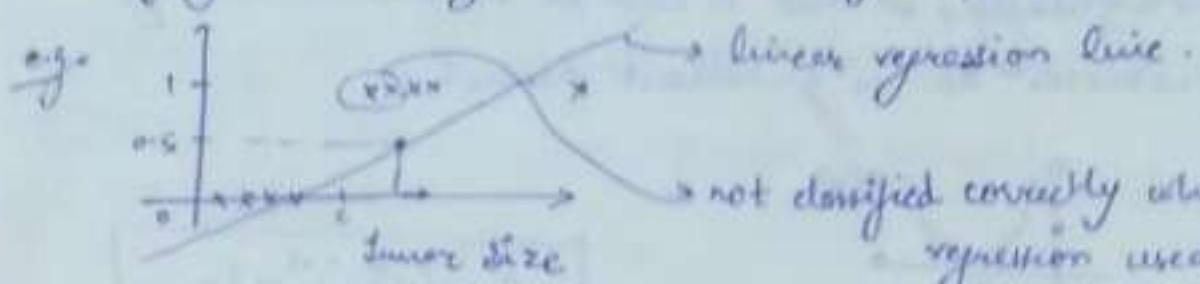
- i) No need to chose  $\alpha$ .
- ii) don't need to iterate.
- iii) is slow when  $n$  is large becoz  $(X^T X)^{-1}$  is  $(n \times n)$  matrix.
- iv) preferred upto atleast  $n = 1000$

If  $X^T X$  is singular or non-invertible matrix, even then the  $\theta$  value can be computed. This condition of non-invertibility generally arises when there are repeated features in  $X$ , or when  $m < n$ , i.e. there are less training examples than there are no. of features. In this case either we delete some features or use regularization which will help fit parameters even when we have a small training set.

→ Logistic Regression:  $0 \leq h_{\theta}(x) \leq 1$

for binary classification problems -  $h_{\theta}(x) > 0.5$ , predict  $y = 1$   
 $< 0.5$ , "  $y = 0$ .

Applying linear regression to a classification problem isn't a good idea.



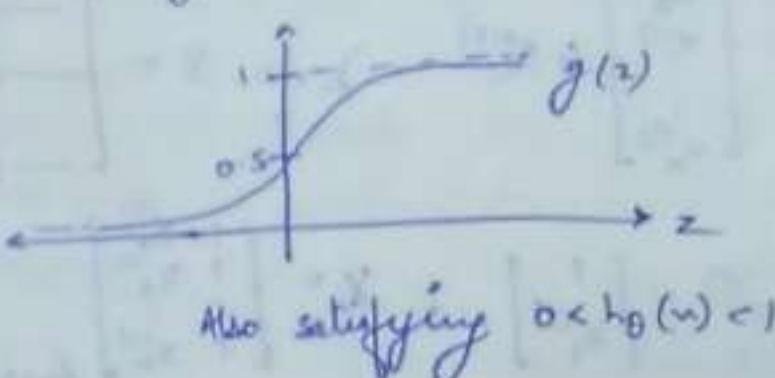
for logistic regression -

$$h_{\theta}(x) = g(\theta^T x)$$

where,  $g(z) = \frac{1}{1+e^{-z}}$  (sigmoid fn/ logistic fn)

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

estimated probability that  
 $y = 1$  on input  $x$ .



e.g. if  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$

$$h_{\theta}(x) = 0.7$$

→ There are 70% chance of tumor being malignant.

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

probability that  $y=1$ , given  $x$ , parameterized by  $\theta$ .

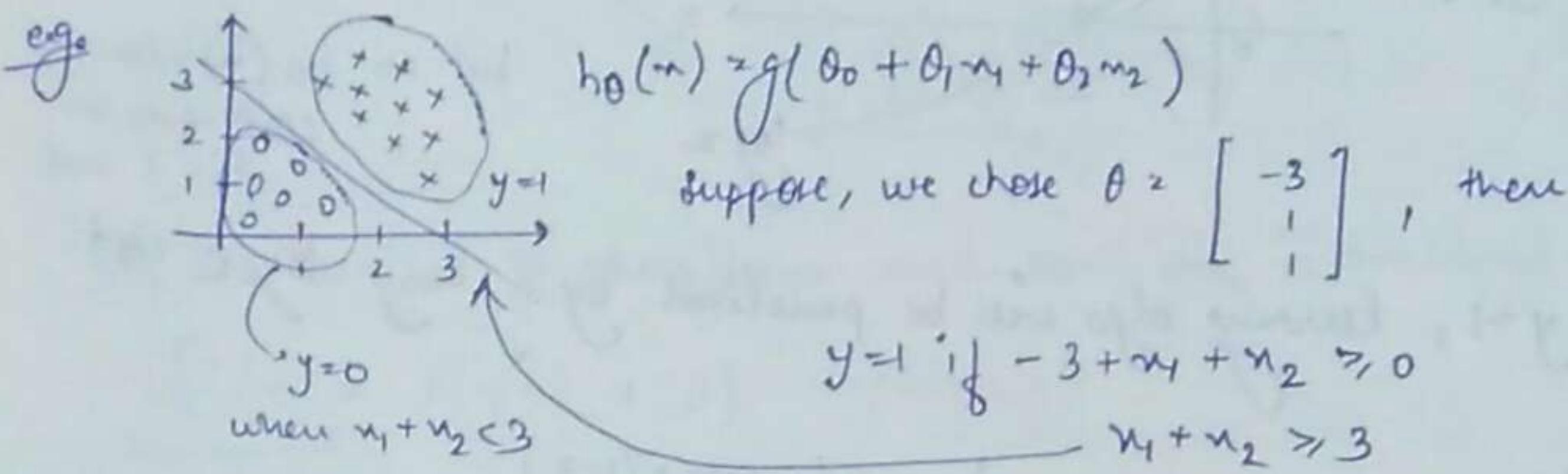
$$y = 0 \text{ or } 1$$

$$\text{so, } P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1$$

when  $x > 0$ ,  $g(x) \geq 0.5$ .

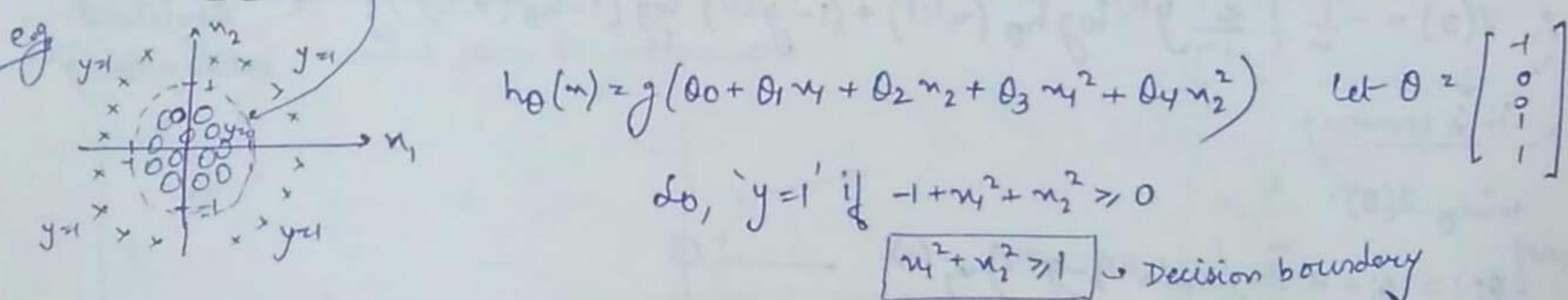
i.e whenever,  $\theta^T x > 0$ ,  $h_\theta(x) = g(\theta^T x) \geq 0.5$

So,  $y=1$  when  $\theta^T x > 0$  and  $y=0$  when  $\theta^T x < 0$ .

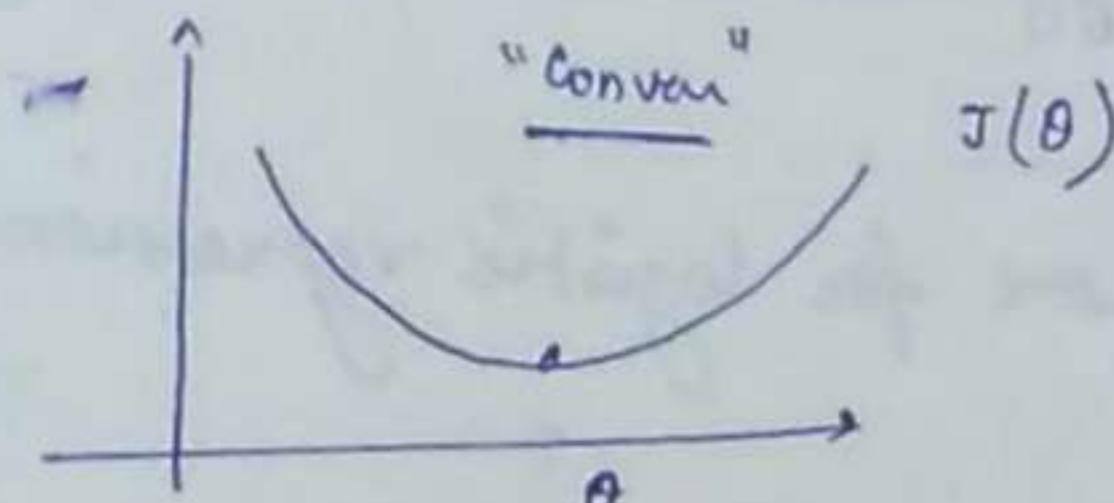
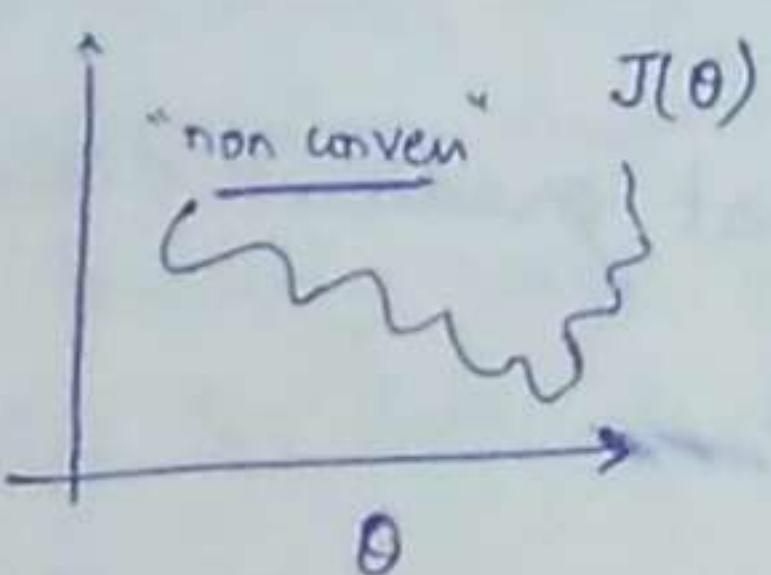


$[n_1 + n_2 = 3] \rightarrow$  'DECISION BOUNDARY' that divides the region of  $y=0$  &  $y=1$ .  
↳ a property of the hypothesis and the parameter  $\theta_0, \theta_1, \dots, \theta_n$  and not a property of the dataset.

Non-linear decision boundaries:



Training set does not define the decision boundary but only used to fit the parameter  $\theta$ , which in turn decides the decision boundary.

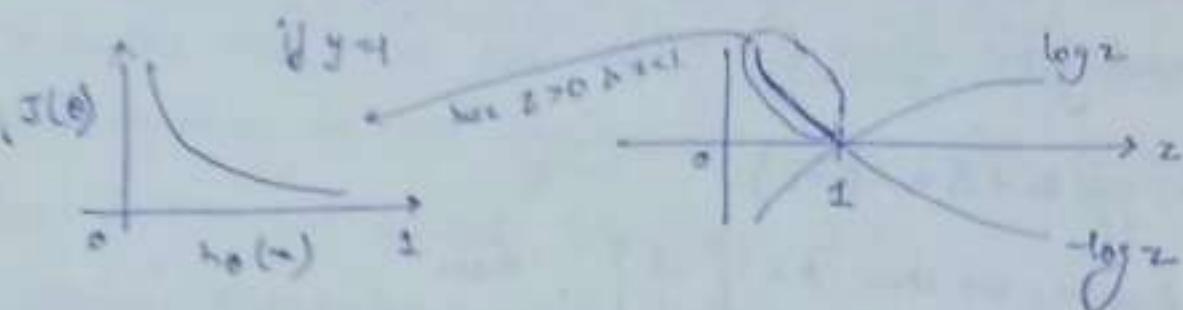


In this case, it is not sure that gradient descent converges to global minimum.

Here, gradient descent will always converge to global minimum.

$J(\theta) = \frac{1}{2} (h_\theta(n) - y)^2$  ends up to be non-convex for  $h_\theta(n) = \frac{1}{1+e^{-\theta^T n}}$  and so, we define new  $J(\theta)$  for logistic regression.

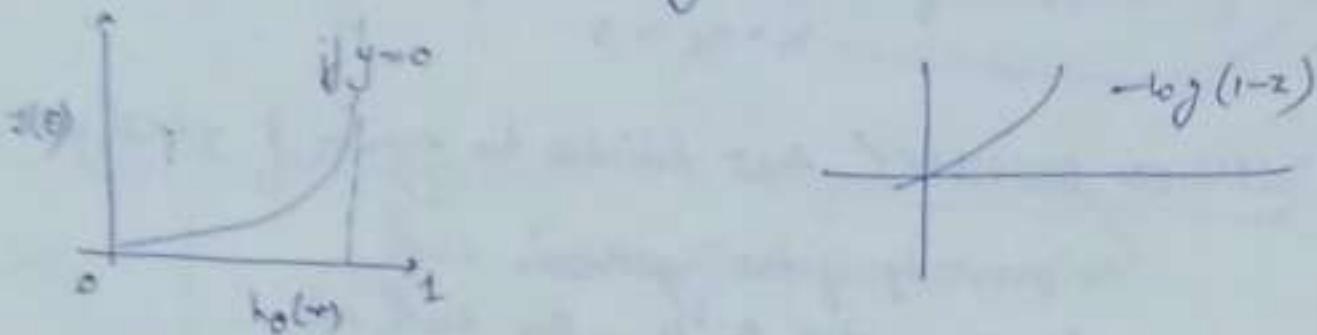
$$J(\theta) = \begin{cases} -\log(h_\theta(z)) & \text{if } y=1 \\ -\log(1-h_\theta(z)) & \text{if } y=0 \end{cases} \quad \text{for single training example.}$$



So, cost = 0 if  $y=1$ ,  
 $h_\theta(z)=1$

but as  $h_\theta(z) \rightarrow 0$   
 $\text{cost} \rightarrow \infty$

If  $h_\theta(z) \approx 0$ , but  $y \approx 1$ , learning algo will be penalized by a very large cost.



$$\boxed{\text{Cost}(h_\theta(z), y) = -y \log(h_\theta(z)) - (1-y) \log(1-h_\theta(z))} \quad \checkmark$$

$$\therefore J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_\theta(z^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(z^{(i)})) \right] \quad \left. \begin{array}{l} \text{minimum} \\ \text{likelihood} \\ \text{estimation} \end{array} \right\}$$

(it is convex)

min <sub>$\theta$</sub>  J( $\theta$ ):

repeat

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_\theta(z^{(i)}) - y^{(i)}) x_j^{(i)} \quad \text{--- (1)}$$

(simultaneously update all  $\theta_j$ )

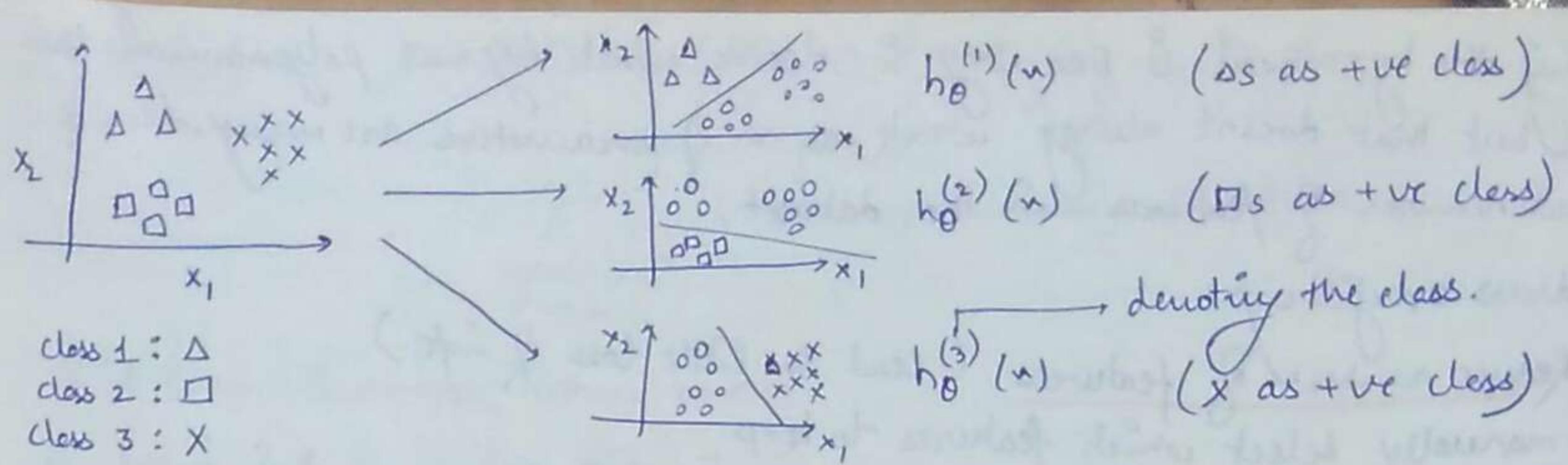
① is same as we got for linear regression also, are linear and logistic same or different algos. the difference is for linear regression  $h_\theta(z) = \theta^T z$  but for logistic  $h_\theta(z) = g(\theta^T z)$ .

feature scaling can be done for logistic regression also so that gradient descent converges fast.

⇒ Optimization algorithm: complex but faster than gradient descent.  $\alpha$  need not be chosen. e.g. conjugate gradient descent, BFGS, L-BFGS.

→ Multiclass classification: One v/s All. or One v/s Rest

for this we turn the problem into three separate 2 class classification problems  
e.g.



So, we have fit 3 classifiers and each one is trained to recognize 1 of 3 classes.

$$h_\theta^{(i)}(n) = P(y=i|n; \theta) \quad (i=1, 2, 3)$$

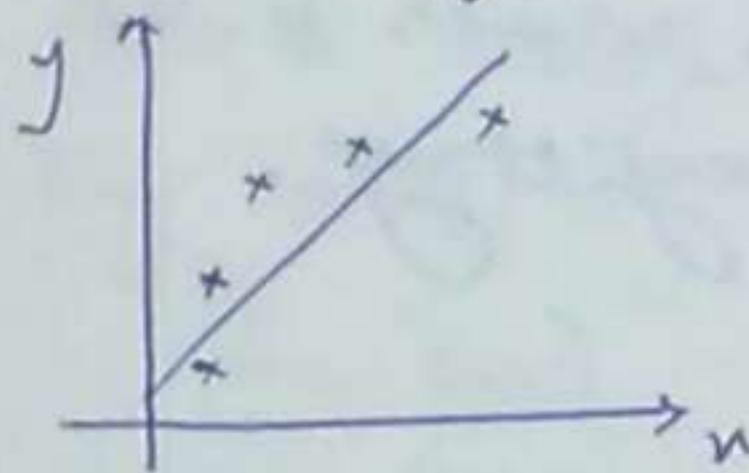
↑  
logistic regression  
classifier

On a new input  $n$ , to make a predict, pick the class  $i$  that minimizes

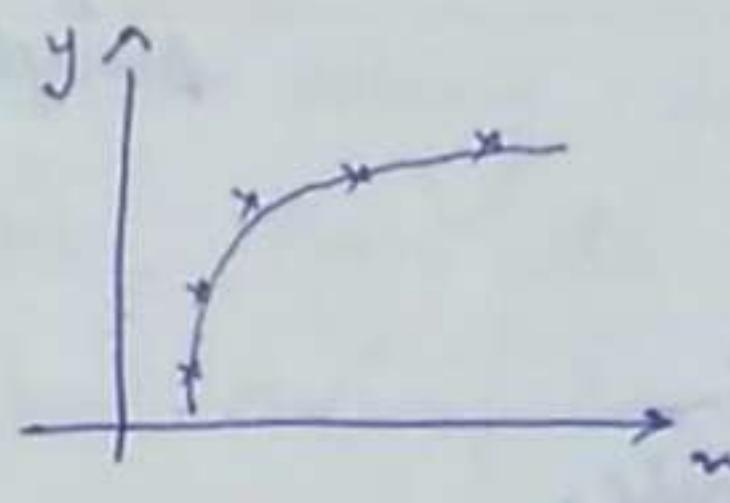
$\max_i h_\theta^{(i)}(n)$ : whichever value of  $i$  gives maximum value of  $h_\theta(n)$  is predicted to be the output with max confidence.

3/5/18

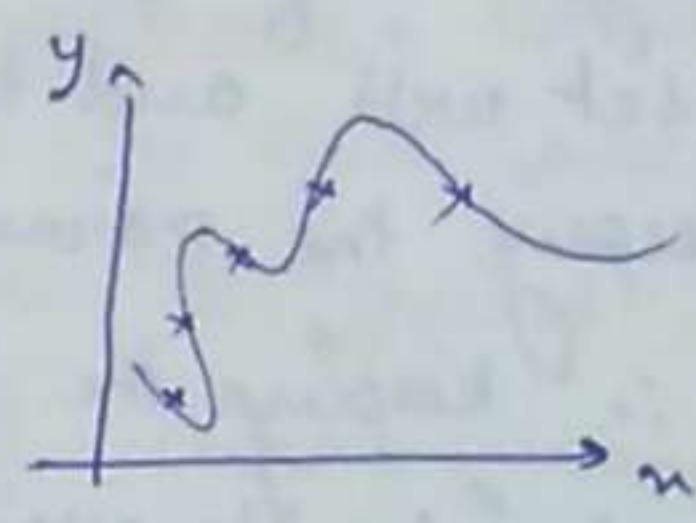
→ Overfitting: causes poor performance.



"Underfit" "high bias"



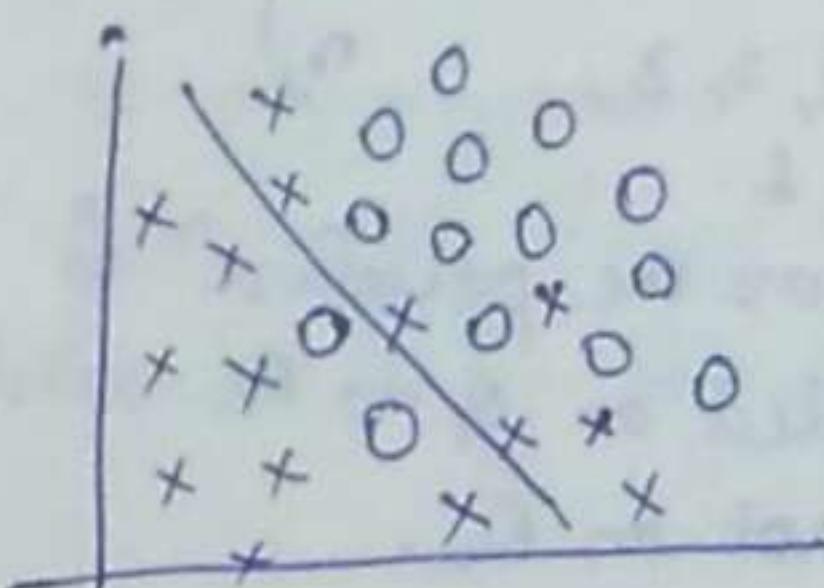
"just right"



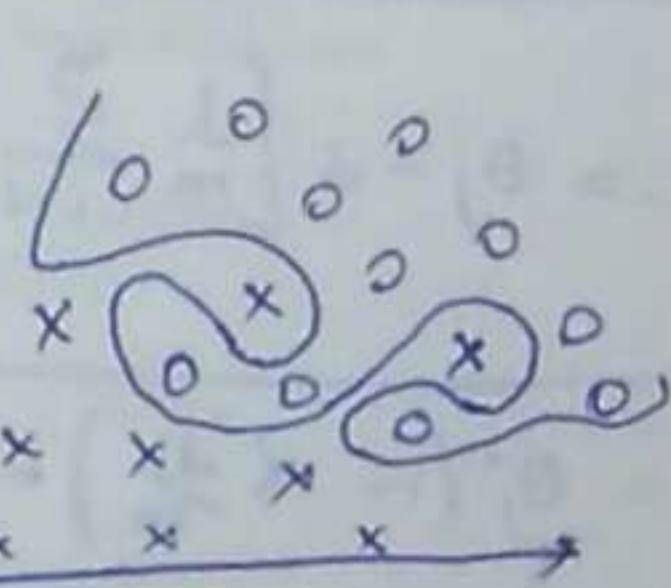
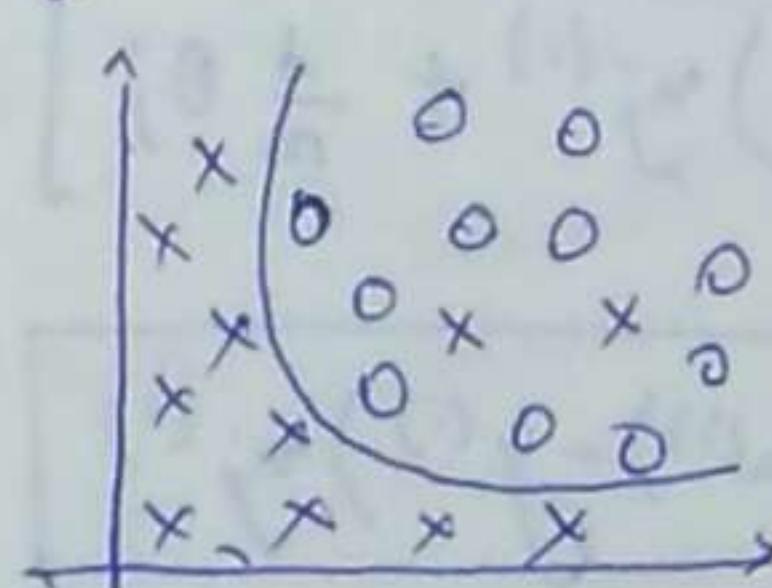
"Overfit" "high variance"

as the model seems to have a bias or preconception that  $y$  is going to linearly vary with  $n$  and despite the data to the contrary.

- happen when we have too many features.  $h(\theta)$  fits training set very well but fails to generalize to new examples.



"Underfit"



in logistic regression e.g.s.

Plotting the hypothesis is one way to decide what degree polynomial to use, but that doesn't always work, as no. of parameters are very close to valid number of features in the dataset.

→ Address overfitting:

1. Reduce number of features (lead to little loss of info.)

- manually select which features to keep.

- model selection algorithm (algs which select features on their own)

2. Regularization

- keep all the features, but reduce magnitude/values of parameters  $\theta_j$ .

- works well when we have lot of features, each of which contributes a bit to predicting  $y$ .

$$\text{for regularize } J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \underbrace{\sum_{j=1}^n \theta_j^2}_{\text{regularization term}} \right]$$

The regularization term helps in controlling the trade off b/w the 2 goals of fitting the training set well and the goal of keeping the parameter small, ∴ keeping the hypothesis relatively simple to avoid overfitting.

$\lambda$  regularization term to shrink all the parameters  $\theta_1$  upto  $\theta_n$ , as we don't know which of them to shrink in order to get simpler hypothesis which is less prone to overfitting.

Note: If  $\lambda$  is set to an extremely large value, algorithm results in underfitting and can't even fit the training set. So,  $\lambda$  should be chosen with care as well. Large  $\lambda$  will lead to setting  $\theta_1, \theta_2, \dots, \theta_n \approx 0$  thus learning  $h(\theta) = \theta_0$ , hence underfitting.

- Regularized linear regression:

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}) w_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = 1, 2, 3, \dots, n)$$

$$\boxed{\theta_j := \theta_j \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}) w_j^{(i)}}$$

but the penalize is applied on  $\theta_1$  to  $\theta_n$  and not on  $\theta_0$ .

$$1 - \alpha \frac{\lambda}{m} < 1$$

(to shrink the parameter a little bit)  
and then do the update as we used to do before

for normal eq:  $\theta = (X^T X + \lambda \begin{bmatrix} 0 & \dots & 0 \\ 0 & \dots & 0 \\ \vdots & \ddots & \vdots \end{bmatrix})^{-1} X^T y$

$\downarrow$   
 $(n+1) \times (n+1)$  dimensional matrix  
 regularization term.

$X^T X$  is non-invertible when  $m < n$ .

If  $\lambda > 0$ ,  $\theta = (X^T X + \lambda \begin{bmatrix} 0 & \dots & 0 \\ 0 & \dots & 0 \\ \vdots & \ddots & \vdots \end{bmatrix})^{-1} X^T y$

$\downarrow$  is invertible.

### - Regularized Logistic Regression:

$$J(\theta) = - \left[ \frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Hence, regularization can take care of overfitting even in case of many features.

$\theta_j$  expression same as that for linear regression with  $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$

## \* Neural Networks: (WEEK 4)

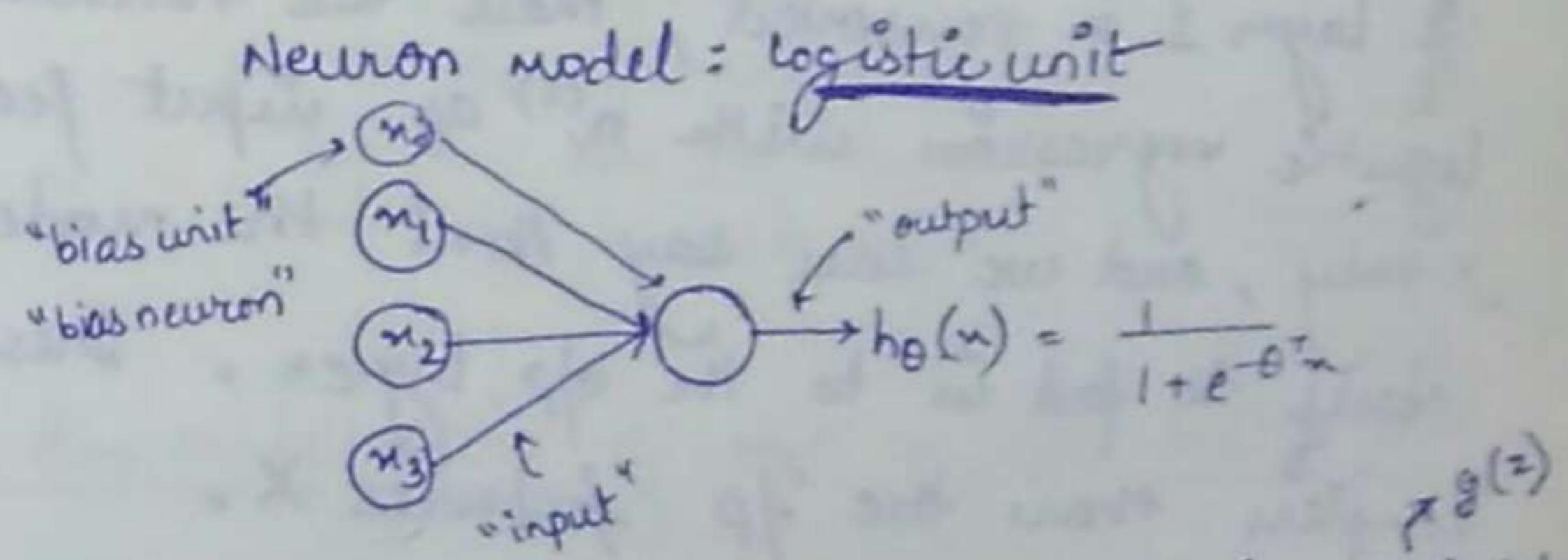
Used when we need to learn complex non-linear hypothesis, when we have a classification problem with many features. In this case if we try apply logistic regression, then its difficult to guess  $\theta_0, \theta_1, \dots, \theta_n$  and not only this but their square and cube terms as well and many others. for  $n=100$ , we will end up around 5000 features to be accommodated in the equation and that too when consider only degree 2 terms.

$x$  or the input vector

for  $20 \times 20$  image is a  $(400 \times 3)$  dimensional vector to accomodate RGB code.

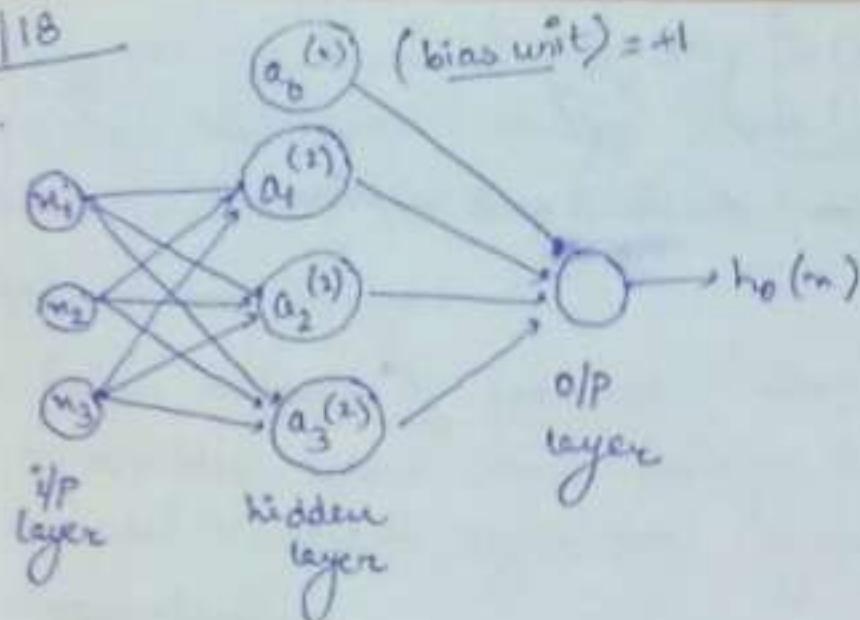
This is too large to be handled by logistic regression.

Neural networks are computationally expensive algorithms.



Artificial neuron is sigmoid or logistic activation function

16/18



$a_i^{(j)}$  = "activation" of unit  $i$  in layer  $j$ .  
 value that's computed  
 by and as o/p by a  
 specific

$\theta^{(j)}$  = matrix of weights controlling  
 for mapping from layer  $j$  to  
 layer  $j+1$ .

$$a_1^{(2)} = g(\theta_{10}^{(1)} n_0 + \theta_{11}^{(1)} n_1 + \theta_{12}^{(1)} n_2 + \theta_{13}^{(1)} n_3) = g(z_1^{(2)})$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} n_0 + \theta_{21}^{(1)} n_1 + \theta_{22}^{(1)} n_2 + \theta_{23}^{(1)} n_3) = g(z_2^{(2)})$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} n_0 + \theta_{31}^{(1)} n_1 + \theta_{32}^{(1)} n_2 + \theta_{33}^{(1)} n_3) = g(z_3^{(2)})$$

$$h_0(n) = a_4^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}) = g(z^{(3)})$$

If network has  $b_j$  units in layer  $j$ ,  $b_{j+1}$  units in layer  $j+1$ , then  $\underline{\theta^{(j)}}$  will be  
 of dimension  $\underline{s_{j+1} \times (s_j + 1)}$ .

Also,  $a^{(1)} = x$

$$z^{(2)} = \theta^{(1)} x = \theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

$$\text{Add } a_0^{(2)} = 1$$

$$z^{(3)} = \theta^{(2)} a^{(2)}$$

$$h_0(n) = a^{(3)} = g(z^{(3)}) \quad [\text{forward propagation}]$$

} eqn of fwd  
 propagation in NN

bcz we forward propagate the activation  
 of i/p layer to calculate the activa? of o/p layer

If layer 1 is removed, then the remaining part of NN is simple logistic regression with  $a^{(2)}$  as input features.  $a^{(2)}$  is learned from  $x$  only, and we can say that the model has learnt the features itself to feed in to the o/p layer. These learnt features are more complex than the i/p features  $x$ .

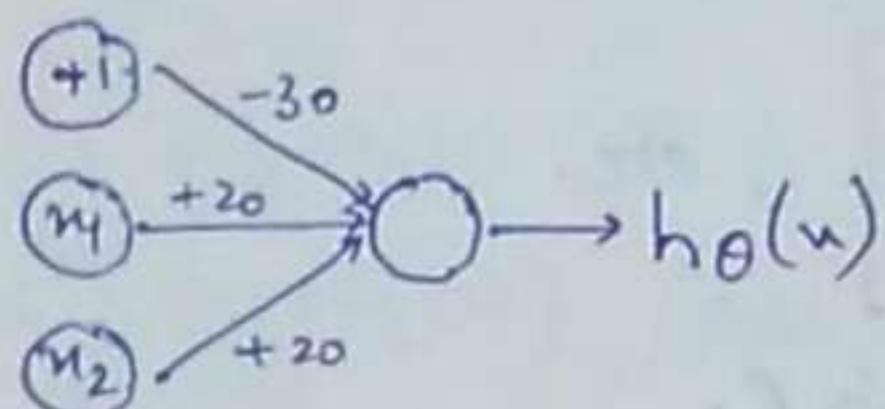
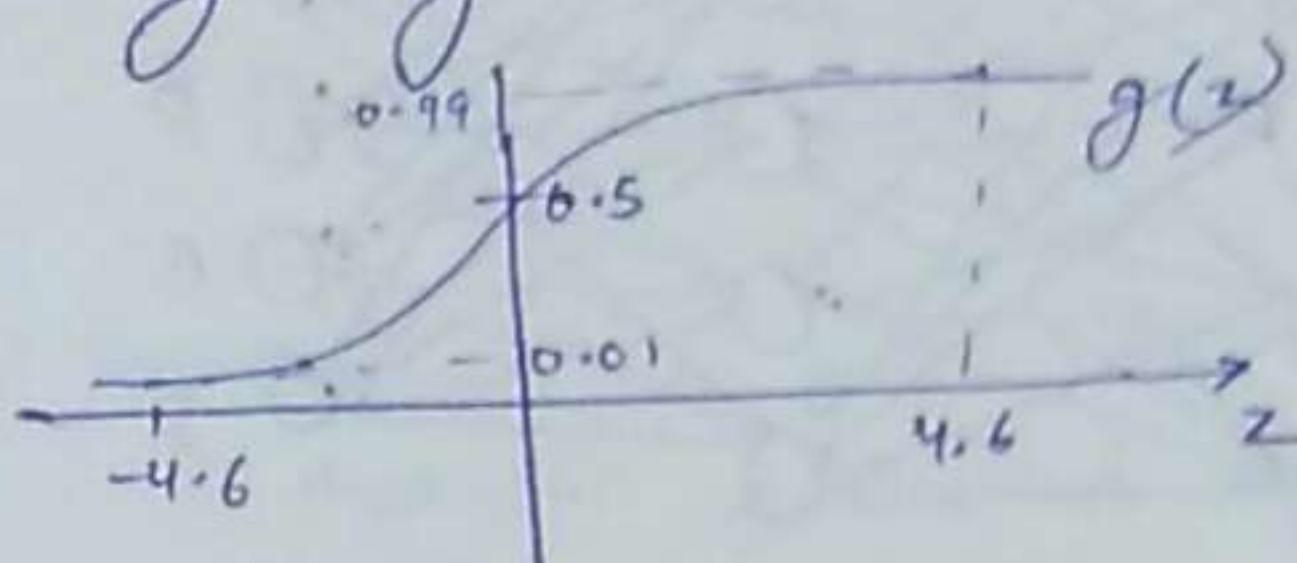
NN architecture : way in which neurons are connected to each other.

A very interesting non linear hypothesis is built by learning more & more complex features layer by layer.

e.g.1 AND

$$y = n_1 \text{ AND } n_2$$

$$n_1, n_2 \in \{0, 1\}$$



$$h_\theta(u) = g(-30 + 20n_1 + 20n_2)$$

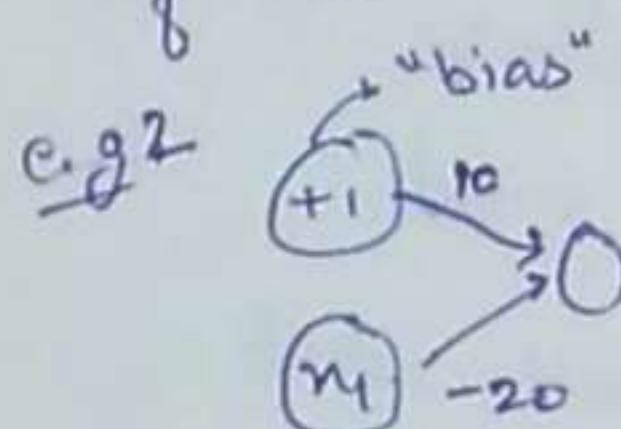
$$\quad \quad \quad \theta_{11} \quad \theta_{12} \quad \theta_{01}$$

$n_1$	$n_2$	$h_\theta(u)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

Hence  $h_\theta(u) \approx X_1 \text{ AND } X_2$

So, the parameter values decide which fn our NN is computing.

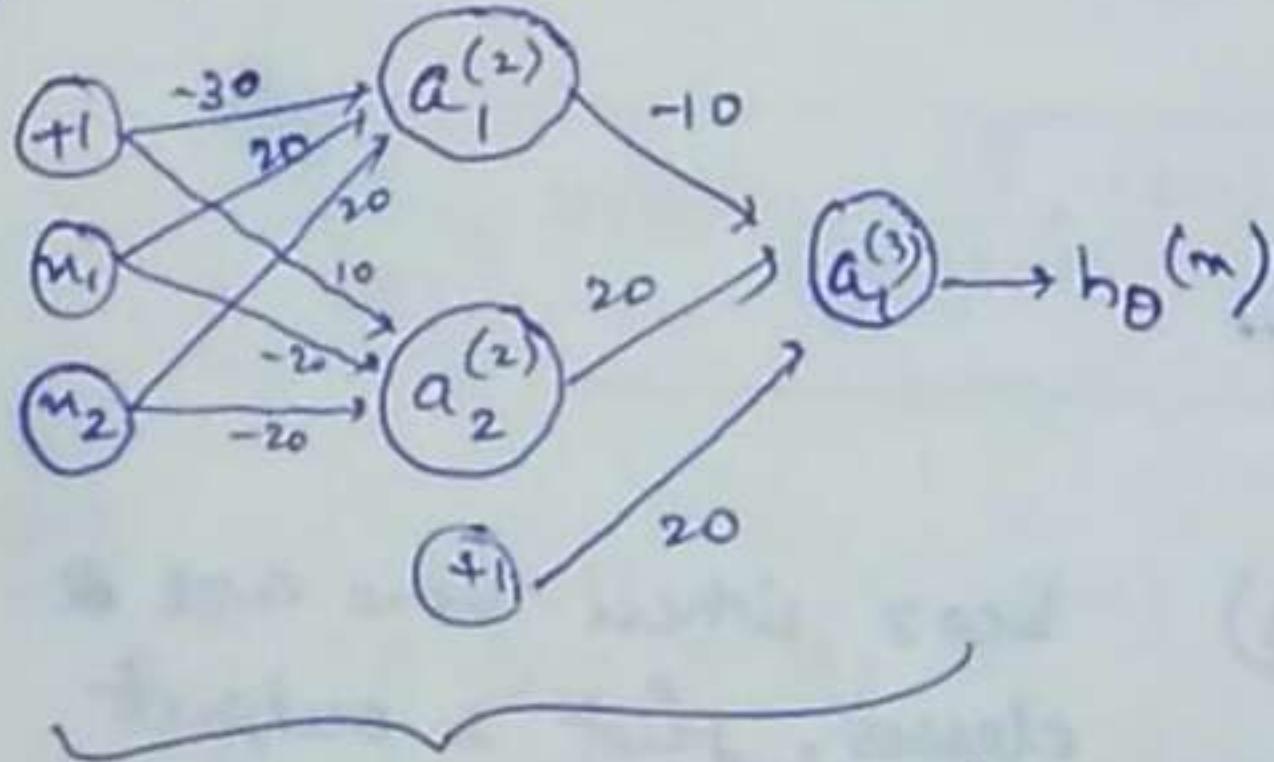
If  $\theta_{10} = -10, \theta_{11} = 20, \theta_{12} = 20, h_\theta(u) \approx X_1 \text{ OR } X_2$



$n_1$	$h_\theta(u)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

Hence  $h_\theta(u) = \text{NOT } X_1$  (Negation of  $X_1$ )

e.g.3  $n_1 \times \text{NOR } n_2$



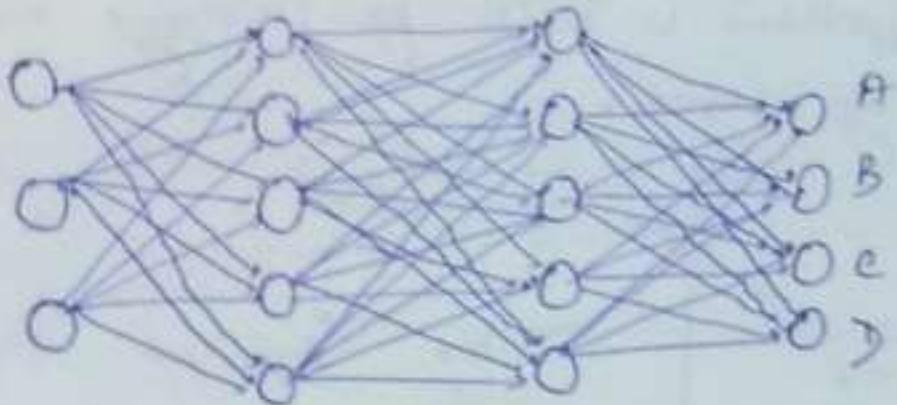
by combining the above e.g. networks, we get this network in e.g.3.

$n_1$	$n_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_\theta(u)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

Slightly more complex fn computed by the NN as compared to  $n_1$  and  $n_2$ .

Multiclass classifier in NN is an extension of one vs all technique.

e.g. to classify if the category is A, B, C or D.



MLP = Multi layer perceptron  
(in MLP, all layers are fully connected)

$$h_\theta(x) \in \mathbb{R}^4 \quad \text{--- (1)}$$

We want  $h_\theta(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ ,  $h_\theta(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ ,  $h_\theta(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ , etc.  
when it's A      when it's B      when it's C.

where training set :  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

$y^{(i)}$  one of  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$  depending on if it's A, B, C or D.

$x^{(i)}$  is the image that is to be classified as A, B, C or D.

$h_\theta(x^{(i)}) \approx y^{(i)}$  where both are  $\mathbb{R}^4$ .

Note  $a_1 + a_2 + a_3 + \dots + a_p \neq 1$  (always)

no. of classes to be identified.

The sum of activation fn need not be 1 bcoz off NN  
are not probabilities.

L = total no. of layers in network,

$s_L$  = no. of units (not counting bias unit) in layer L.

Binary classification :  $y = 0 \text{ or } 1$ .  $s_L = 1$

Multiclass " :  $y = K$  possible values  $s_L = K$  ( $K \geq 3$ ) bcoz when there are 2 classes, just 1 output layer may work.

for neural network,  $J(\theta)$  is generalized of the logistic regression :-

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1-y_k^{(i)}) \log(1-(h_\theta(x^{(i)}))_k) \right] + \frac{1}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_L} \sum_{j=1}^{s_{L+1}} (\theta_{j,i}^{(l)})^2$$

where  $(h_\theta(x))_i = i^{\text{th}}$  output  
and  $h_\theta(x) \in \mathbb{R}^K$

Backpropagation algorithm tries to minimize this cost function.  
To minimize  $J(\theta)$ , we need code to compute  $J(\theta)$  and  $\frac{\partial}{\partial \theta_{ij}} J(\theta)$  where  $\theta_{ij} \in \mathbb{R}$  are parameters.

Let's use figure ① for one training example  $(x, y)$  and write eq's of back propagation through vectorized implementation:

$$\begin{aligned} a^{(1)} &= x \\ z^{(2)} &= \theta^{(1)} a^{(1)} \\ a^{(2)} &= g(z^{(2)}) \quad (\text{add } a_0^{(2)}) \\ z^{(3)} &= \theta^{(2)} a^{(2)} \\ a^{(3)} &= g(z^{(3)}) \quad (\text{add } a_0^{(3)}) \\ z^{(4)} &= \theta^{(3)} a^{(3)} \\ a^{(4)} &= h_\theta(x) = g(z^{(4)}) \end{aligned}$$

Gradient computation: back propagation algo

$\delta_j^{(l)}$  = "error" of node  $j$  in layer  $l$ .

for each o/p unit (layer  $L = 4$ )

$$\boxed{\delta_j^{(4)} = a_j^{(4)} - y_j} \Rightarrow \underbrace{\delta^{(4)}}_{\text{vectors where each has dimension equal to no. of o/p units in our network.}} = a^{(4)} - y$$

called back propagation  
bcz 1st  $\delta^{(4)}$  is calculated,  
then using it  $\delta^{(3)}$  is  
and so on. bcz we are  
proceeding from back, it  
is called so.

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \star g'(z^{(3)})$$

element wise multiplication

$$\text{where } g'(z^{(3)}) = a^{(3)} \star (1 - a^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} \star g'(z^{(2)})$$

No  $\delta^{(1)}$  term as 1st layer is ip layer.

$$\boxed{\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)}}$$

ignoring 1 term of regularization

2/6/18

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$$

{overall error for all training examples}

$$\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

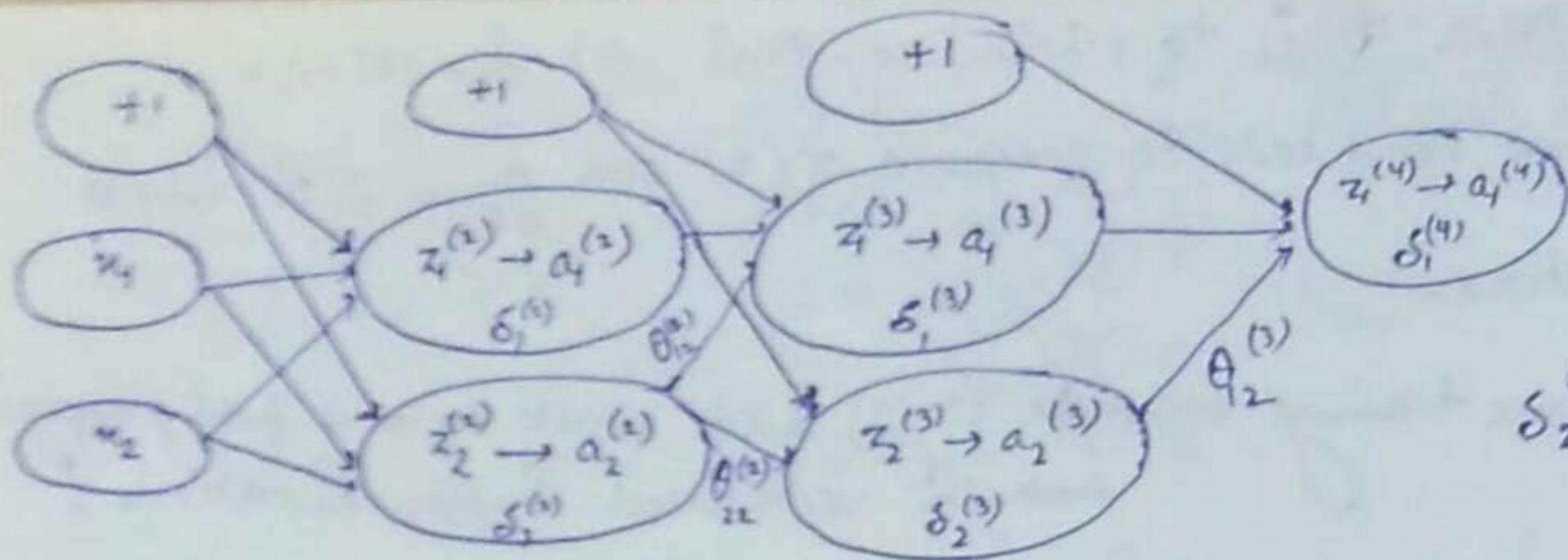
{vectorized implementation}

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \quad \text{if } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0$$

$$\text{and } \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$$

Backpropagation is somewhat mathematically a black box only.

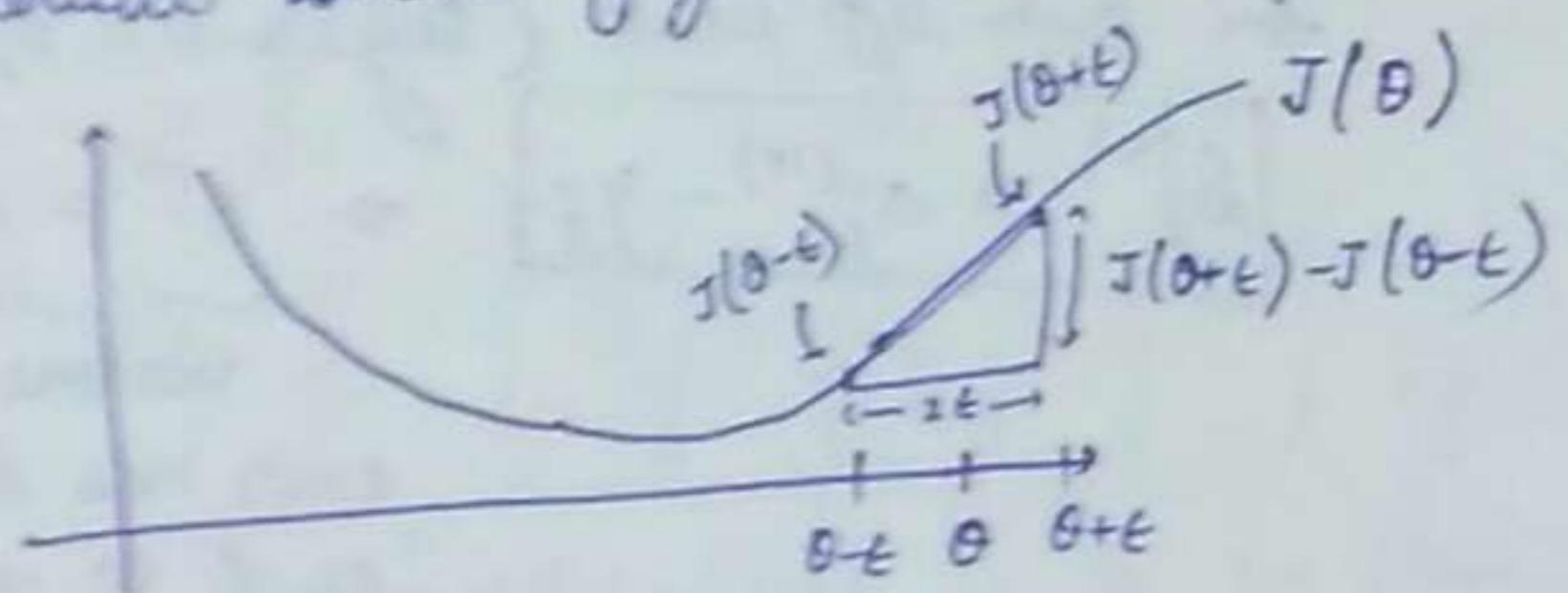


$$\delta_2^{(3)} = \theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

and similarly the others

- Gradient checking: to make sure that the implementation of fwd prop and back prop is 100% correct.

Numerical estimate of gradients -



$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

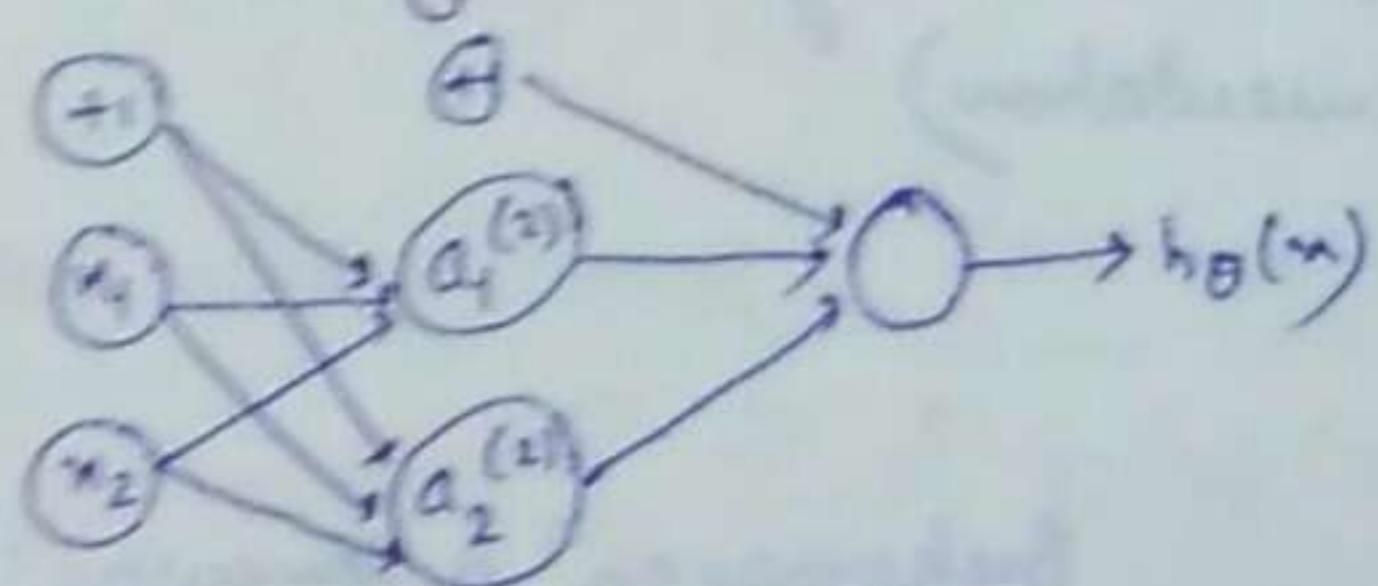
if  $\theta$  is parameter vector, then  $\frac{\partial}{\partial \theta_i} J(\theta) \approx \frac{J(\theta_i + \epsilon, \theta_1 \dots \theta_n) - J(\theta_i - \epsilon, \theta_1 \dots \theta_n)}{2\epsilon}$

Imp: check if  $\frac{\partial}{\partial \theta} J(\theta) = \text{DVec}$  as calculated from backprop.  
(unrolled  $D^{(1)}$ ,  $D^{(2)}$ ,  $D^{(3)}$ )

this gradient checking needs to be done once and then should be disabled before training the classifier, becz if gradient iter is run on every compute, code will be very slow.

Initializing  $\theta$  as zeros(n,1) in NN doesn't work as it works in case of logistic regression.

In case of zero initialize:



we find that if  $\theta_{ij}^{(l)} = 0$  initially if i,j,l then;

$$a_1^{(2)} = a_2^{(2)}. \text{ Also } \delta_1^{(1)} = \delta_2^{(2)}$$

$$\frac{\partial}{\partial \theta_{01}^{(1)}} J(\theta) = -\frac{\partial}{\partial \theta_{02}^{(1)}} J(\theta)$$

This way everytime for every iteration the weights coming from  $\Theta_1$  will remain the same, weights coming from  $\Theta_2$  will remain the same and so will those coming from  $\Theta_3$ , i.e. all the hidden units are computing exactly the same features leading to a highly redundant representation.

Hence, we initialize the parameters of NN with random nos, this random initialize is also called symmetry breaking.

Initialize each  $\theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$

3/6/18

↳ different from the & previously introduced.

No. of i/p units : dimension of features  $n^{(1)}$

No. of o/p units : No. of classes

Reasonable default : 1 hidden layer, and if more than 1 then same no. of hidden units in every layer. More no hidden units, better it is.

### Training a neural network:

- 1) Randomly initialize weights.
- 2) Implement feed propaga to get  $h_\theta(x^{(i)})$  for any  $x^{(i)}$ . {i.e o/p for any given i/p}
- 3) Implement code to compute cost fn  $J(\theta)$
- 4) Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \theta_j^{(l)}} J(\theta)$

A for loop is used to implement backprop on all examples atleast for the first time of usage. Other alternatives for 'for' loop are complex.

- 5) Use gradient checking to compare  $\frac{\partial}{\partial \theta_j^{(l)}} J(\theta)$  computed using backprop vs using numerical estimate of gradient  $\frac{\partial}{\partial \theta_j^{(l)}} J(\theta)$ . Then disable gradient checking code as its computationally very slow.
- 6) Use Gradient descent or advanced optimiza<sup>n</sup> method with backprop to try to minimize  $J(\theta)$  as a fn of parameters  $\theta$ .

Note for NN,  $J(\theta)$  is non convex and optimiza<sup>n</sup> algos get stuck at local minima but it is found that these algos generally find the best local minima, if not the global minimum.

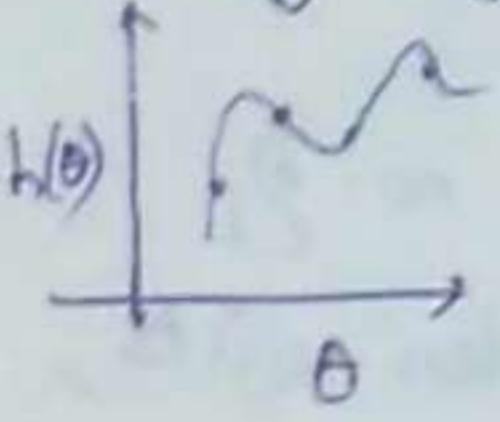
## \* WEEK 6 :

Suppose, an algo applied doesn't give good results on test set, then;

- 1) Get more training examples OR
- 2) Try smaller set of features (to avoid overfitting) OR
- 3) Try getting additional features (\* " underfitting) OR
- 4) Try adding polynomial features. ( $w_1^2, w_2^2, w_1w_2, \dots$ )
- 5) Try  $\lambda$  big or  $\lambda$  very small.

There are easy techniques to select which of the above options one should take to improve the model, they are ML diagnostics.

Evaluating a hypothesis i.e evaluating the hypothesis learned by the learning algorithm.

 Plotting a hypothesis and then analyzing it is possible only when it's in 2D but when we have many parameters, this evaluation technique fails.

So, what we do is split the data into train and test set. We then learn the parameters from the training set and then calculate  $J_{\text{test}}(\theta)$ .

Misclassified error :  $\text{err}(h_\theta(x), y) = \begin{cases} 1 & \text{if } h_\theta(x) > 0.5, y = 0 \\ 0 & \text{or if } h_\theta(x) < 0.5, y = 1 \\ \text{otherwise} & \end{cases}$  error (0/1 misclassified error)

$$\text{Test error} = \underbrace{\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_\theta(x_{\text{test}}^{(i)}), y^{(i)})}_{\text{This denotes the fraction of the test set examples that have been mislabeled by our hypothesis.}}$$

Model selec<sup>n</sup> problems: deciding value of  $\lambda$ , deciding which features to keep or which to reject.

- Training set error is not a good indicator of how well the hypothesis will do on new examples, bcoz the training set error is likely to be lower than the actual generalization error.

- Suppose we want to decide which degree of polynomial should be selected ( $d$ ). For this, we find  $J_{\text{test}}(\theta)$  for different values of  $d$  and select that value of  $d$  which gives minimum  $J_{\text{test}}(\theta)$ . Hence, this is fitting the parameter ' $d$ ' on test data. So, there is less chances that this error is similar to the generalized error i.e. with the examples that the model has not seen before. Hence, to solve this issue, we divide our data into train set, valida<sup>n</sup>/cross valida<sup>n</sup> set, test set. where approx 60% is train, 20% is valida<sup>n</sup> & 20% is test data.

Now,

$$\text{training error : } J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\text{cross valida<sup>n</sup> error : } J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

$$\text{Test error : } J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

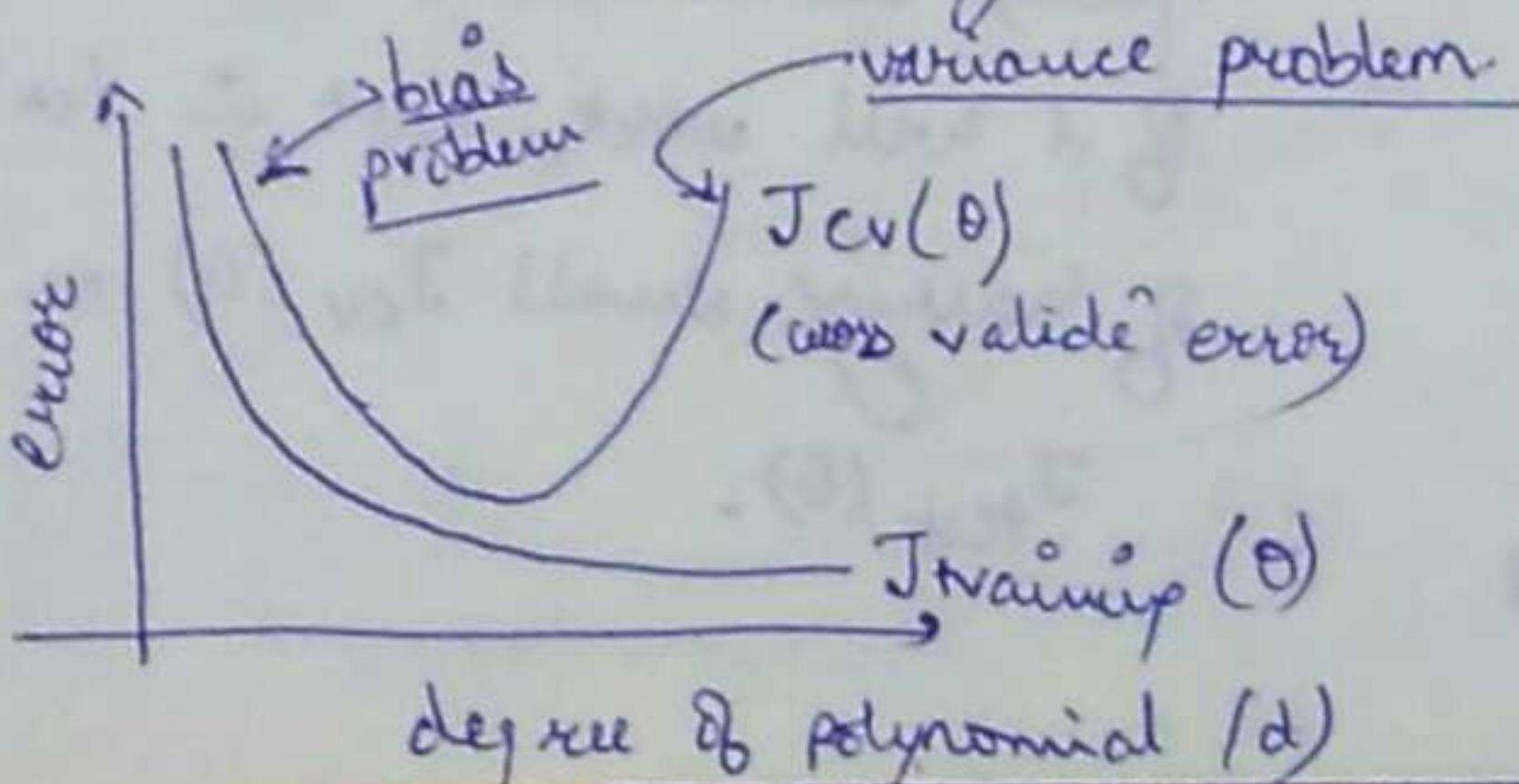
So, now instead of selected using test set for model selec<sup>n</sup>, we will use cv set for model selection. Now the  $d$  is fitted on cv set and so now test set can be used to estimate the generaliza<sup>n</sup> error of the model that was selected.

Hence, it is a wrong practice to use the test set to select the model and then use the same value to generalize the model. It is always advised to divide the data into 3 parts.

- High bias / High variance:

(underfit)      (overfit)

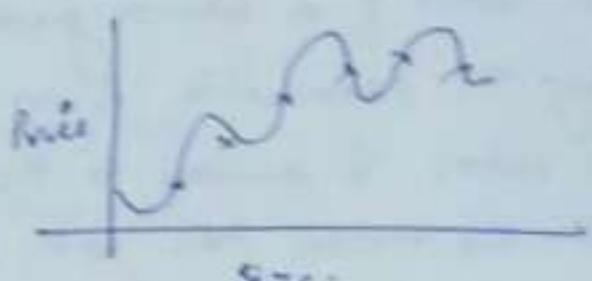
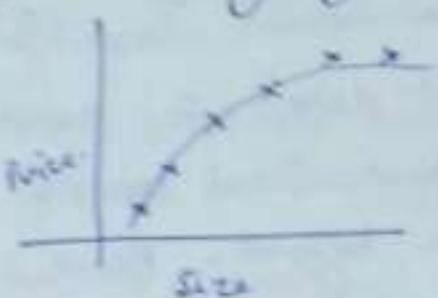
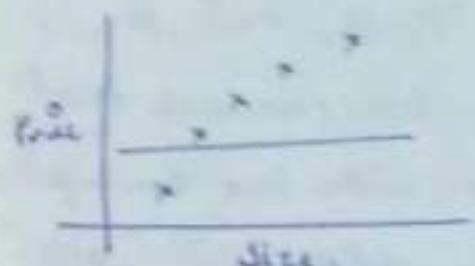
Q Suppose our learning algo is performing less well than we were hoping. (i.e  $J_{\text{cv}}(\theta)$  or  $J_{\text{test}}(\theta)$  is high). Is it a bias or a variance problem?



Bias (underfit) :  $J_{\text{train}}(\theta)$  as well as  $J_{\text{cv}}(\theta)$  will be high.

Variance(overfit) :  $J_{\text{train}}(\theta)$  is low while  $J_{\text{cv}}(\theta)$  will be high.

Regularize helps prevent overfitting but how does it affect the bias & variances of a learning algorithm.



Large  $\lambda$   
High bias (underfit)  
 $\lambda = 10000, \theta_0 \approx 0, \theta_1 \approx 0, \dots$   
 $h_{\theta}(x) \approx \theta_0$

Intermediate  $\lambda$   
"Just right"

Small  $\lambda$   
High variance (overfit)  
 $\rightarrow \lambda \approx 0$

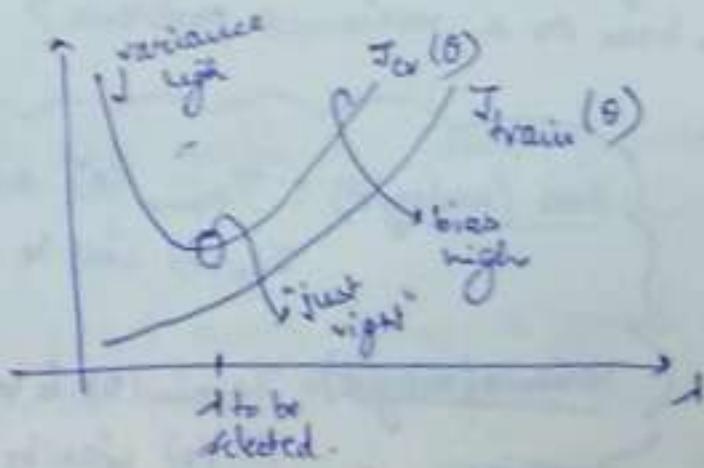
$$\text{given: } h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2}_{\text{regularize term}}$$

Choosing the regularize parameter ' $\lambda$ ':

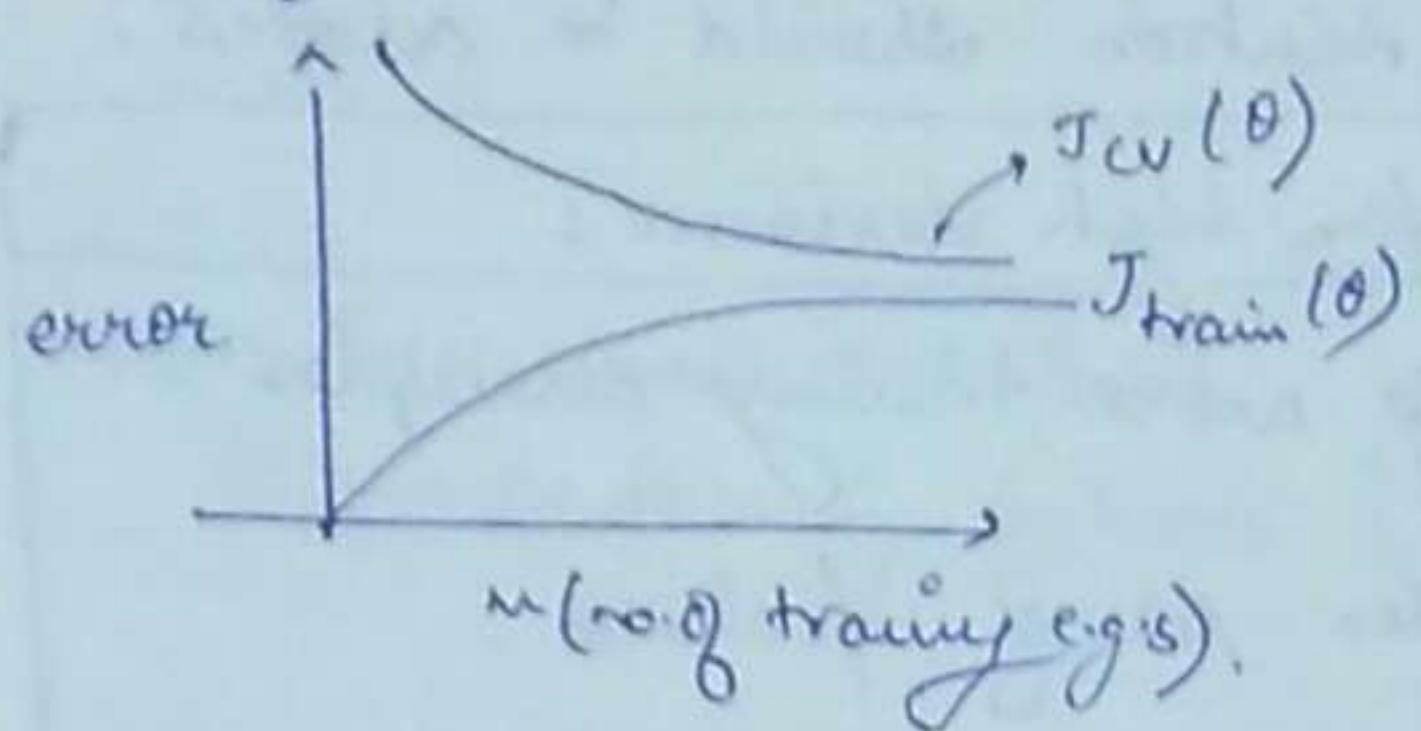
- 1. Try  $\lambda = 0$
  - 2. Try  $\lambda = 0.01$
  - 3. Try  $\lambda = 0.02$
  - 12. Try  $\lambda = 10$
- } Similarly trying on a set of values  $\lambda$  gives 12 different models

Now, every  $\lambda$  will be corresponding to a parameter vector  $\theta$  by minimizing  $J(\theta)$ . Now these parameter vectors corresponding to the different  $\lambda$  values are used to compute  $J_{CV}(\theta)$  and the one giving minimum value for  $J_{CV}(\theta)$  is selected. Now this  $\theta$  vector is used to compute  $J_{test}(\theta)$  to report the error. So, that's model selection applied to selecting the regularize parameter ' $\lambda$ '.

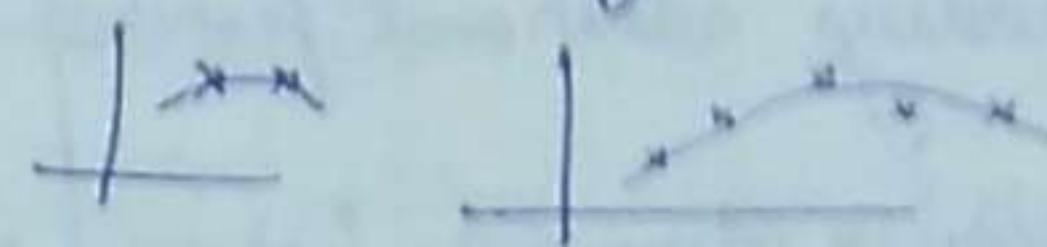


Hence, some intermediate value of  $\lambda$  will work best in terms of having small  $J_{CV}(\theta)$  or small  $J_{test}(\theta)$ .

→ Learning curves : to check and improve the performance of algs.



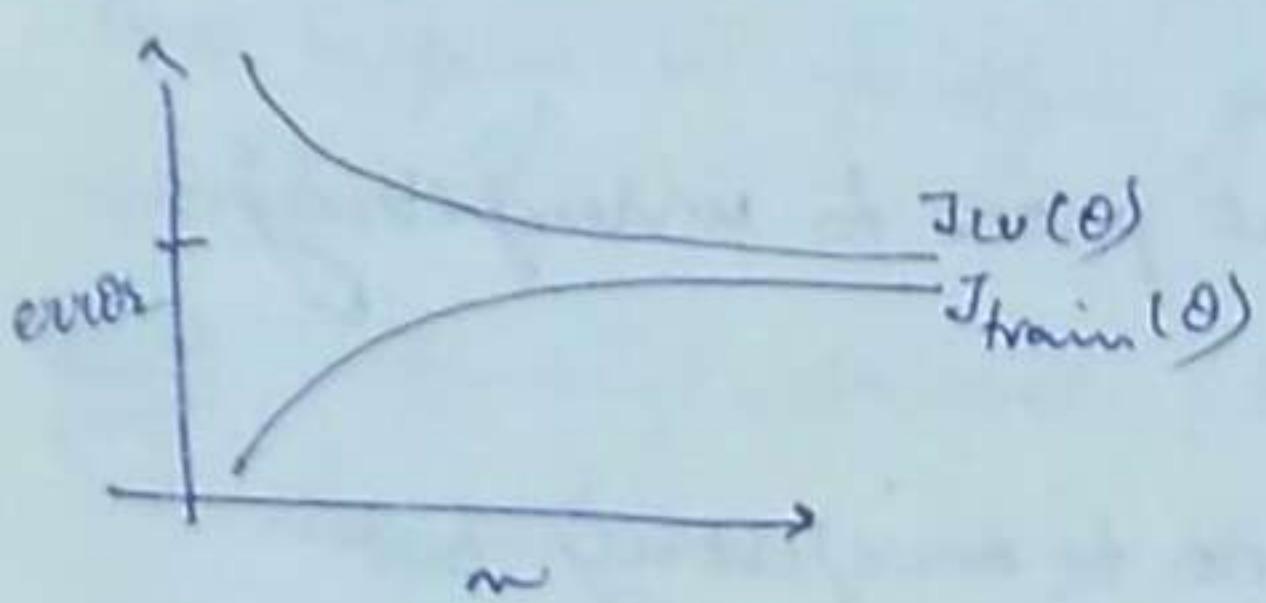
This is bcoz when we have less training ex's it is easier to fit them in an eq.



but when n is large, it gets harder to fit the hypothesis to them.

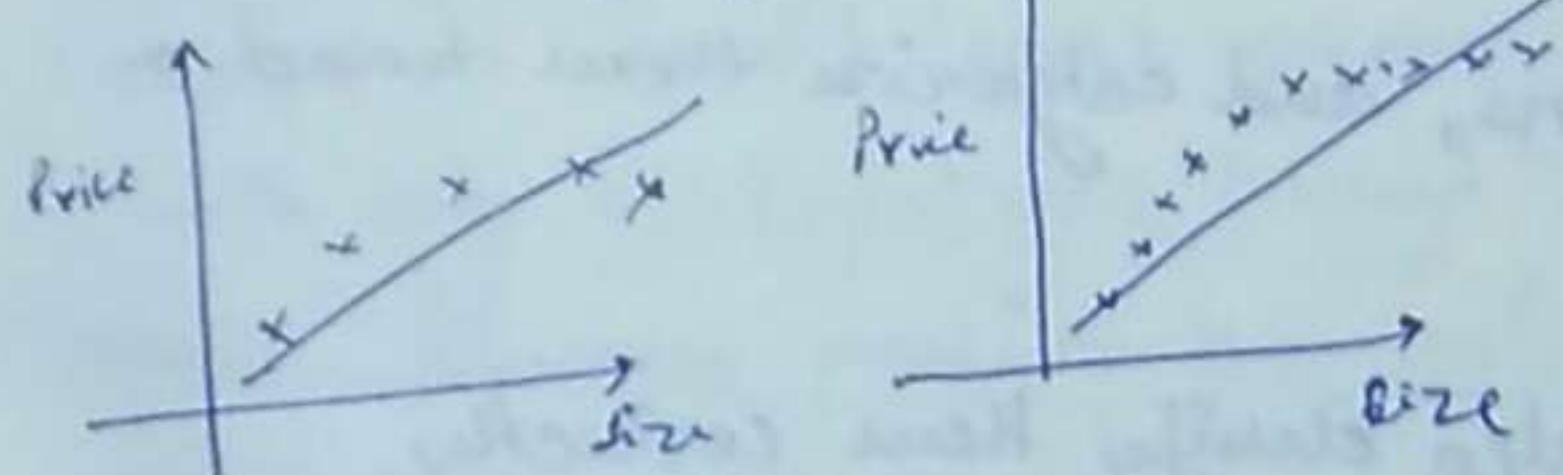
$J_{CV}(\theta)$  fits with tiny n bcoz with greater n, a <sup>more</sup> general hypothesis can be made.

### High bias case



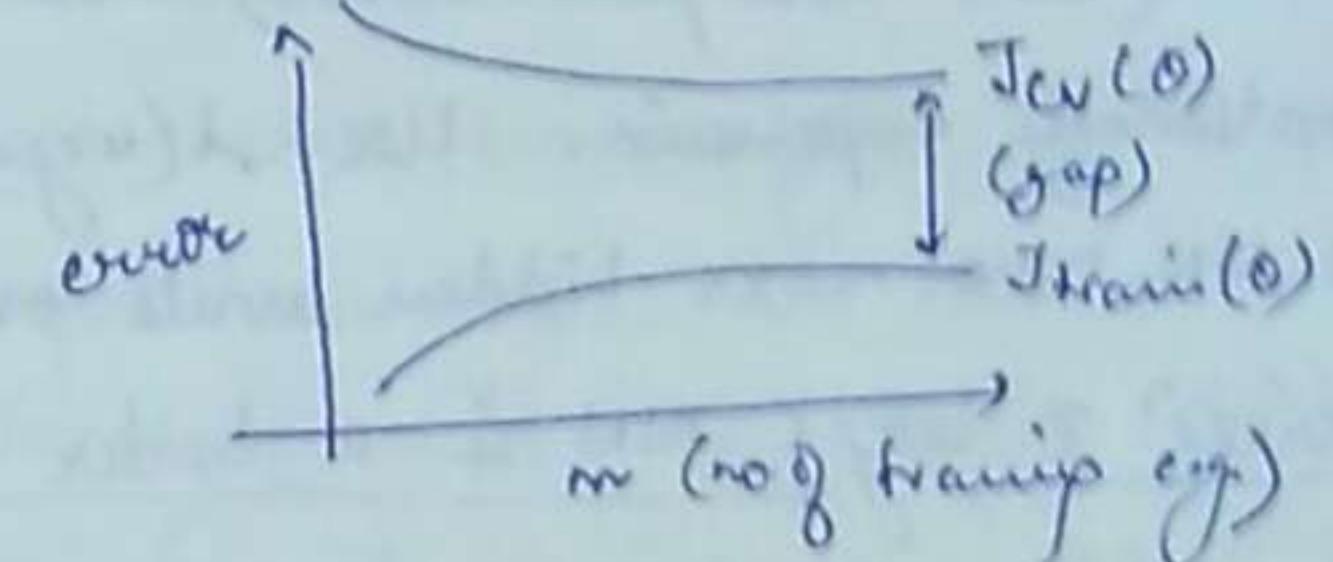
In case of high bias, performance of the hypothesis will be quite similar on the training set as well as the testing i.e. cross validate set.

In this case both  $J_{CV}(\theta)$  and  $J_{train}(\theta)$  are high and close to each other. Hence if a learning alg is suffering from high bias, getting more training examples won't help much bcoz as we can see  $J_{CV}(\theta)$  is not tiny much with tiny n, it is only getting flattened bcoz in case of high bias:

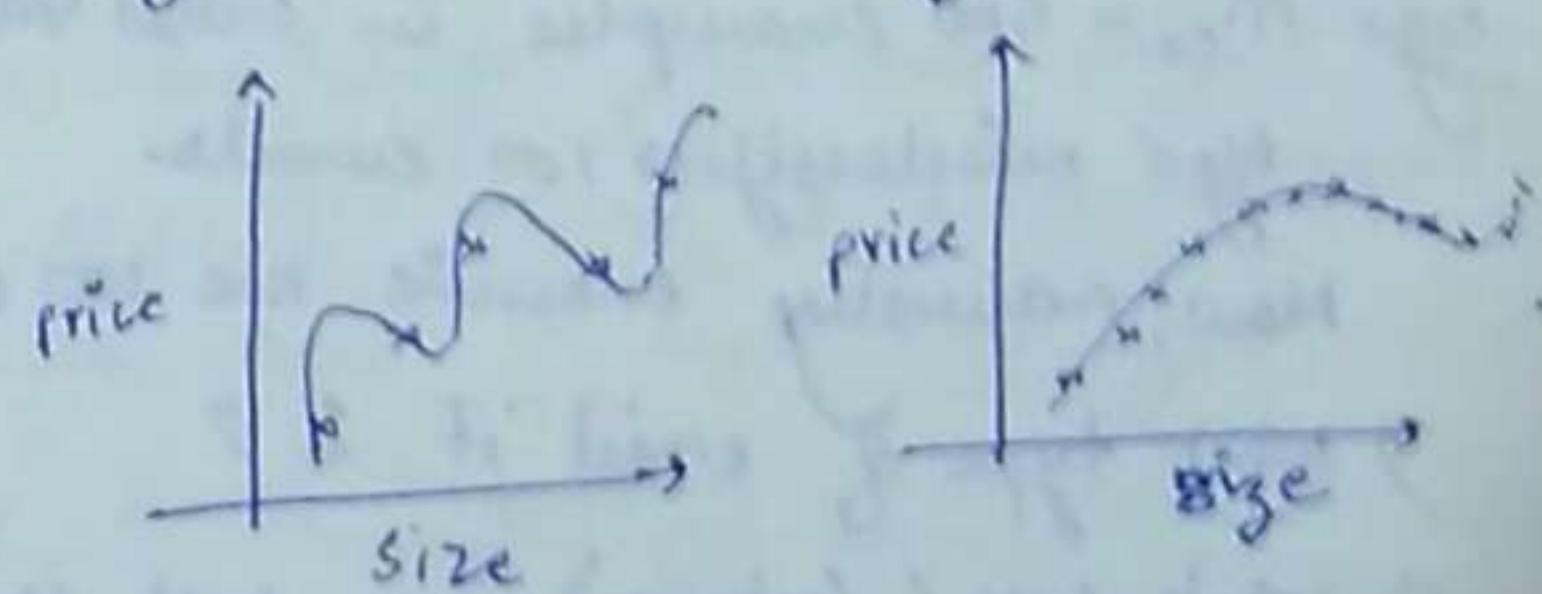


$$h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

### High variance case



In case of high variance, getting more training data helps, bcoz  $J_{CV}(\theta)$  and  $J_{train}(\theta)$  seems to be converging towards each other with tiny n. This is bcoz in case of high variance  $h_{\theta}(x)$  seems to less overfit the data in case of greater no. of training examples.



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

(small t)

Hence, drawing these learning curves help us know if there is a bias or a variance problem and hence what solution should be adopted.

### \* So to fix high bias:

- Try adding additional features
- Try adding polynomial features
- Try  $\lambda = 1$ .

### \* To fix high variance:

- Try getting more training examples
- Try smaller set of features.
- Try  $\lambda = 1$ .

## → Neural networks and overfitting:

- Small NN (having fewer parameters) are more prone to underfitting but computationally cheaper.
- Large NN (with more parameters) are more prone to overfitting but computationally expensive. Use  $\lambda$  (regularizer) to address overfitting. Large NN have more hidden units or more hidden layers. Using regularizer in large NN is a better choice than using small NN.

## → Error analysis: (to be done on CV set & not test set)

Premature optimize in computer programming can be dealt with by help of plotting learning curves to know if the data is suffering from high bias or high variance in advance, rather than going by gut feeling.

Error analysis refers to manually examine the examples (in cross validate set) that our algo made errors on and finding out the systematic trend if any of examples the algo makes errors on.

e.g.  $M_{cv} = 500$  examples in cross validate set  
Algo misclassifies 100 emails.

Now, manually examine the 100 errors, and categorize them based on:

- 1) What type of email it is?
- 2) What cues/features can help the algo classify them correctly.

Suppose, emails are found to be

Phishing : 12 Replica/fake : 4 intended to steal password : 53 Other : 31	
	Deliberate misspellings : 5 Unusual email writing : 16 Unusual use of punctuations : 32

Suppose these emails are found to have the following features

Now, an algo can be designed that focuses on identifying the unusual punctuation use and help in better identifying the spam mails, thus helping in fruitful guidance towards making a successful model.

Note "Stemming software (Porter stemmer)" helps treating discount/discounts/discounting/discounted as same words. But it might hurt when it considers universe and university to be same.

To check if stemming is doing well for our algo, we can calculate the error with and without stemming and then easily decide whether to involve stemming or not, and similarly for other kind of s/w and ideas as well, bcoz checking manually everything would be tough.

Note This initial exercise of error analysis and all is a powerful tool for deciding where to spend the time next bcoz first we look at the errors it makes and do error analysis to see what other mistakes it makes and use that to inspire further development.

Skewed classes : when examples of 1 class highly outnumber the examples of other class. In such cases using just one single no. evaluation doesn't work properly. e.g. if there is a dataset with only 0.5% patients having cancer and I have an algo which always predicts the patient to be healthy, then this single no. evalution gives me algo is 99.5% correct which is actually not the case.

Hence in case of skewed classes, this single value evaluation fails.  
So, we come up with Precision / Recall evaluation metric.

		① Actual class ②	
		True +ve False +ve	
Predicted class	①	True +ve	False -ve
	②	False -ve	True -ve

y=1 in case of rare class that we want to detect.

Precision: Of all the patients where we predicted  $y=1$ , what fraction actually has cancer?

$$\frac{\text{True positives}}{\# \text{ predicted } +ve} = \boxed{\frac{\text{True } +ve}{\text{True } +ve + \text{ False } +ve}}$$

\* High precision is good.

sensitivity: true +ve rate  
specificity: true -ve rate

Recall: of all the patients that actually have cancer, what fraction did we correctly detect as cancer?  
or  
Sensitivity

$$\boxed{\frac{\text{True positives}}{\text{True } +ve + \text{ False } -ve}} = \frac{\text{True positives}}{\# \text{ actual positives}}$$

\* High recall is good and specificity =  $\boxed{\frac{\text{True negatives}}{\# \text{ actual negative}}}$   
So, precision/recall will give a better sense of how well our classifier is working. A classifier with high precision and high recall is actually a good classifier as it covers the whole scene and can't cheat by predicting  $y=0$  or  $y=1$  all the time. But if it predicts  $y=0$  all the time, then precision↑ but recall is low i.e. a bad classifier. If  $y=1$  all the time then high recall but low precision.

→ Trading off precision & recall: (tweak the threshold value)

Suppose we want to predict  $y=1$  (cancer) only if very confident.

One way to do that is to modify the hypothesis  $h_0(n)$  from

$$h_0(n) \geq 0.5 \text{ Predict 1}$$

$$h_0(n) < 0.5 \text{ Predict 0}$$

$$\text{to } h_0(n) \geq 0.7 \text{ Predict 1}$$

$$h_0(n) < 0.7 \text{ Predict 0} \quad i.e.$$

we say patient has cancer only if he has 70% or more chance of having cancer, leading to higher precision but in contrast this classifier will have lower recall bcz now  $y=1$  is predicted for a smaller no. of patients.

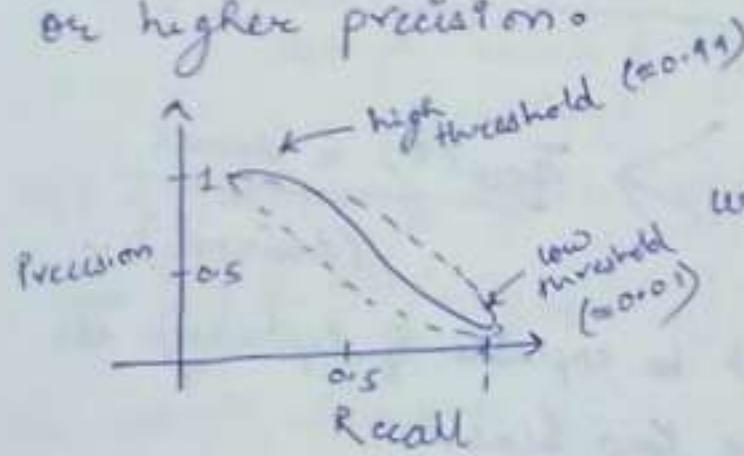
But, suppose we want to avoid missing too many cases of cancer (i.e. avoid false -ve) bcz suppose a patient has cancer & we fail to tell it to them, it can be really bad as then they might not get the desired treatment. So, in case of doubt also we want to predict  $y=1$ , in this case we change  $h_0(n)$  to

$$h_0(n) \geq 0.3, \text{ Predict 1}$$

$$h_0(n) < 0.3, \text{ Predict 0}$$

In this case, we have higher recall but lower precision.

Hence, it depends from case to case if we need higher recall or higher precision.



where  $h_0(n) \geq \text{threshold}$ .

This precision-recall graph can take many shapes depending on the classifier details.

- Q. How to decide this threshold automatically and how to decide which algo is better like in case below:

	Precision	Recall	Avg.
Alg 1	0.5	0.4	0.45
Alg 2	0.7	0.1	0.4
Alg 3	0.02	1.0	0.51

\* \* In case of single no. exist this problem didn't arise but now we have two nos. to decide upon.

To convert this precision-recall to a single value eval metric:  
 we can calculate  $\boxed{\text{Average} = \frac{\text{Precision} + \text{Recall}}{2}}$  and the one with  
 higher avg. is the better algorithm. But this always doesn't work,  
 well as seen in previous example. Algo 3 has such low precision  
 but still is assumed to be best based on avg. value.

So, we come up with  
 to bring the reality  
 out of precision &  
 recall

$$\boxed{\text{F Score or } F_{\text{Score}} = \frac{2PR}{P+R}}$$

If either P or R is zero F score is 0, so  
 for a good F score both P and R need to be  
 reasonably high. Perfect F score = 1  
 when P = R = 1.

### ④ How much data to train on?

"It's not who has the best algo that wins, but who has the most data".

Note: Using a learning algo with many parameters help having  
 low bias and then using a large dataset helps having low  
 variance.

\* low bias  $\Rightarrow T_{\text{train}}(0)$  is small  $\rightarrow T_{\text{test}}(0)$  is small  
 \* low variance  $\Rightarrow T_{\text{train}}(0) \approx T_{\text{test}}(0)$  (desired)

Note parameters of the learning algo must be capable of capturing all  
 the features of the dataset to have low bias.

Hence, large such parameters with large dataset give high accuracy  
algs.

$\rightarrow$  This can be done by comparing it with human, i.e. given features n,  
 can a human successfully predict y. If yes, that means the  
 features are or the parameters are sufficient for good accuracy alg.

Note: logistic and linear regression along w/ NN are low bias algos.

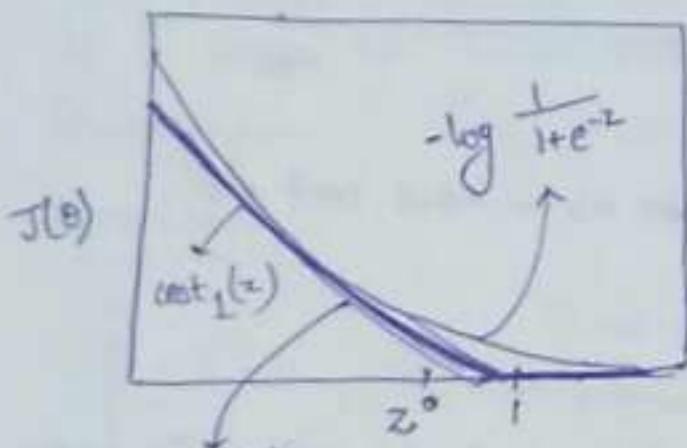
## \* WEEK 7:

4/6/18

### → Support Vector Machines: (Large Margin Classifier)

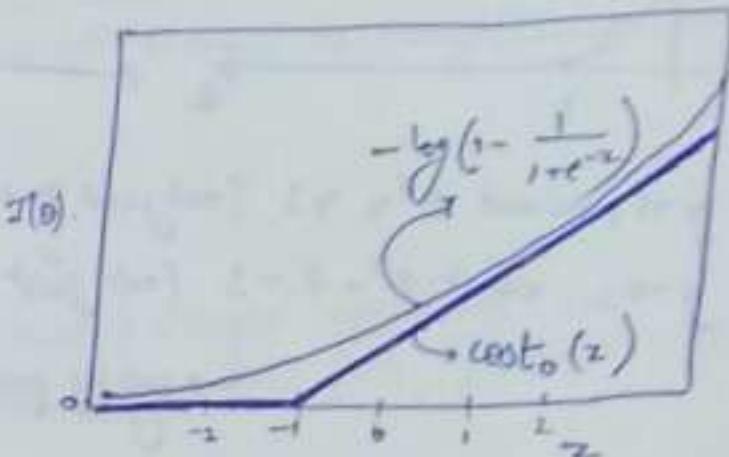
This algo gives a more cleaner and powerful way of learning complex linear functions. It can be obtained by modifying logistic regression a bit.

if  $y=1$  (want  $\theta^T x \geq z > 0$ ):



straight line (new cost fn when  $y=1$ , pretty similar to the cost fn in case of logistic regression)

if  $y=0$  (want  $\theta^T x < 0$ ):



$cost_1(z)$ : cost corresponding to  $y=1$ .  
 $cost_0(z)$ : " . . . . . "  $y=0$ .

(cost fn in case of logistic regression was:

$$\min_{\theta} \frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1-y^{(i)}) (\log(1-h_{\theta}(x^{(i)}))) \right] + \frac{1}{2m} \sum_{j=0}^n \theta_j^2$$

So, cost fn for SVM is:

$$\min_{\theta} \frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} cost_1(\theta^T x^{(i)}) + (1-y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=0}^n \theta_j^2 \quad \text{--- (1)}$$

$$= \min_{\theta} C \sum_{i=1}^m [y^{(i)} cost_1(\theta^T x^{(i)}) + (1-y^{(i)}) cost_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=0}^n \theta_j^2 \quad \text{(for SVM)} \quad \text{--- (2)}$$

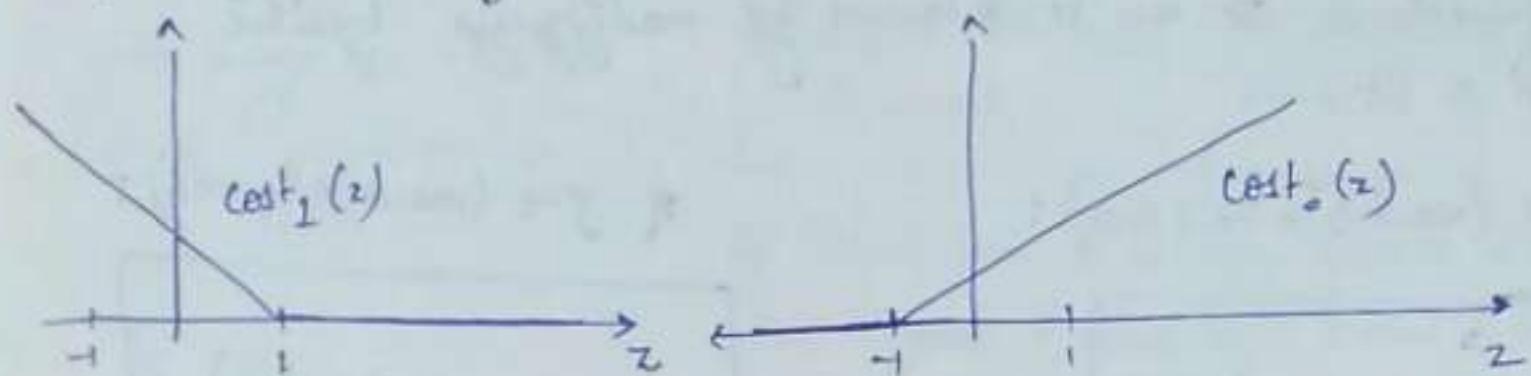
instead of notation  $A+B$ , it is written as  $C+A+B$  where if  $C$  is less  $B$  is giving more weight just as when  $A$  was more  $B$  was given more  $C$  is also regularization parameter (work like  $\lambda$ )

weight. Also  $m$  is removed becz it was in multiplica<sup>t</sup> and removing it does not change the desired minimum value of  $\theta$ .

If we put  $c = 1$ , the cost fn ① and ② turn out to be same.

So,  $h_0(x)$  for SVM is  $\begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$  unlike that of logistic

regression where  $y = 1$  if  $h_0(x) \geq 0.5$  and  $y = 0$  otherwise.



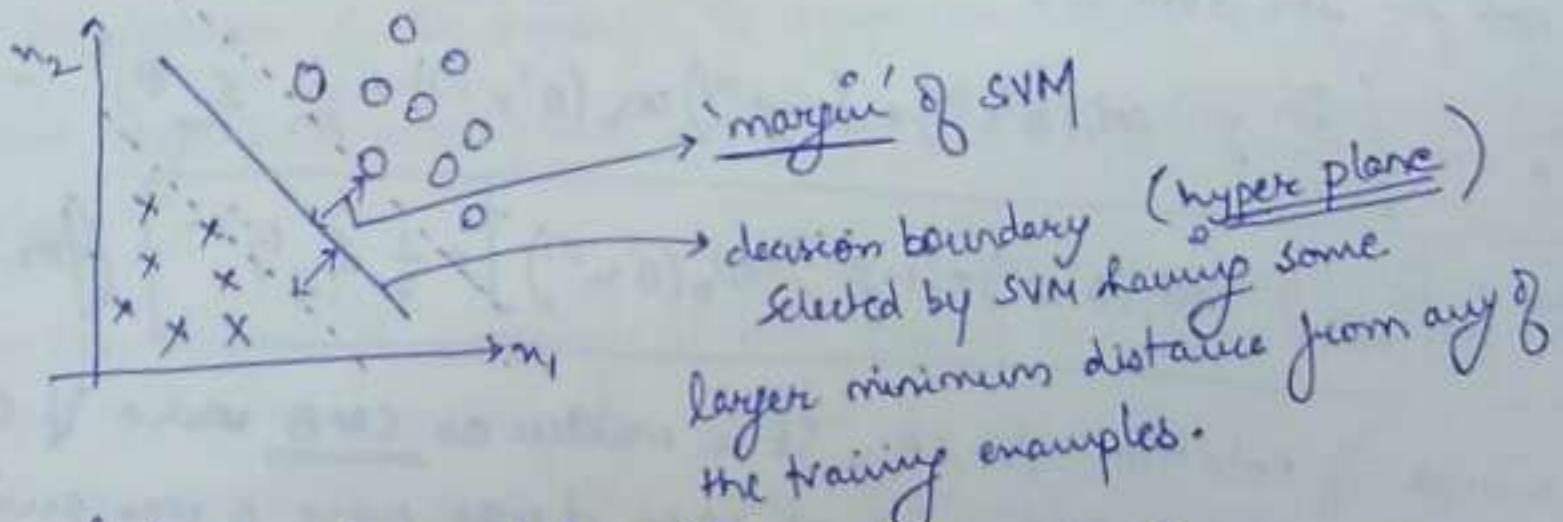
if  $y=1$ , want  $\theta^T x \geq 1$  (not just  $\geq 0$ )  
if  $y=0$ , want  $\theta^T x \leq -1$  (not just  $\leq 0$ )  
do as to have cost = 0.  
safety purpose.

If  $c$  is chosen to be very very large then the minimize<sup>t</sup> objective will choose a value so that  $A = 0$ .

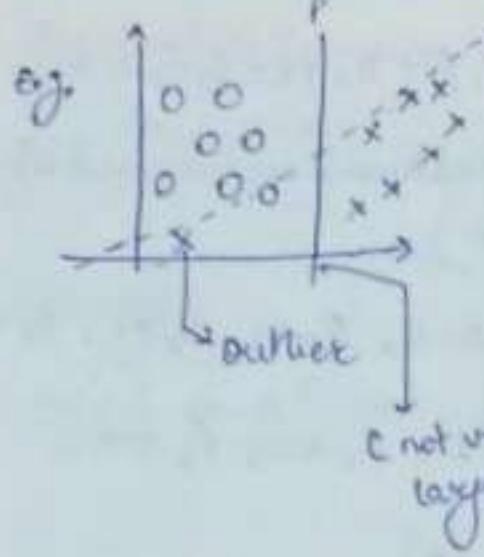
$\Rightarrow \min_{\theta} \frac{1}{2} \sum_{i=1}^m \theta_j^2$  such that  $\theta^T x^{(i)} \geq 1$  if  $y^{(i)} = 1$  and  $\theta^T x^{(i)} \leq -1$  if  $y^{(i)} = 0$ .

On solving this optimize<sup>t</sup> problem, when minimizing this as fn of parameter  $\theta$ , we get a very interesting decision boundary.

e.g. SVM decision boundary : linearly separable case



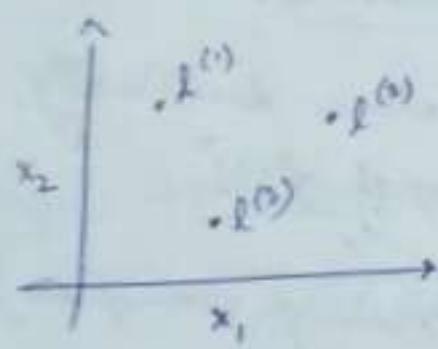
SVM tries to separate the data with as large margin as possible.



Using a large margin classifier may be quite sensitive to outliers. If  $c$  is very large (i.e.  $\lambda$  is small) ' $q$ ' boundary will be selected which is actually not a good choice to make just becaz of one outlier. This decision boundary is actually a sign of overfitting. ' $p$ ' is the right decision boundary and it will be selected when the regularization parameter  $c$  is reasonably small i.e.  $\lambda$  is large (avoiding overfit).

Even if the data is not linearly separable, SVM does its work well.  
It is better to ignore outliers than to change to overfitting decision boundary.

### - Kernels:



Given  $n$ , compute new features depending on proximity to landmarks  $l^{(1)}, l^{(2)}, l^{(3)}$ .

$$\text{Given } n: f_1 = \text{similarity}(n, l^{(1)}) \\ = \exp\left(-\frac{\|n - l^{(1)}\|^2}{2\sigma^2}\right)$$

where  $\|n - l^{(1)}\|$  denotes Euclidean distance b/w the two.

$$f_2 = \text{similarity}(n, l^{(2)})$$

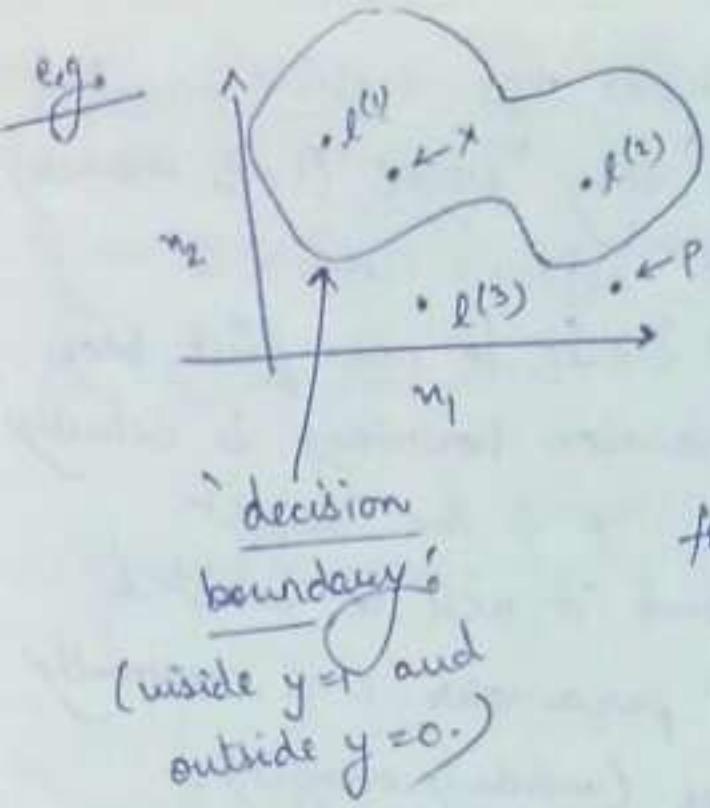
$$f_3 = \text{similarity}(n, l^{(3)}) = K(x, l^{(3)})$$

"Kernel function" (Gaussian Kernel)

$$\text{if } n \approx l^{(1)}: f_1 \approx \exp\left(-\frac{\sigma^2}{2\sigma^2}\right) \approx 1$$

$$\text{if } n \text{ is far from } l^{(1)}: f_1 \approx \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

Each of the landmarks define a new feature, i.e. a given training e.g. we can compute three features  $f_1, f_2, f_3$  given  $l^{(1)}, l^{(2)}, l^{(3)}$ .



e.g. Predict 1 when  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 > 0$

(where  $f_1, f_2, f_3$  are features corresponding to  $l^{(1)}, l^{(2)}$  and  $l^{(3)}$ )

Suppose we calculated  $\theta_0 = 0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$   
for the training e.g.  $X$ ;  $f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$

$$\text{So, } \theta_0 + \theta_1(1) + \theta_2(0) + \theta_3(0) \\ = 0.5 + 1 = 0.5 > 0$$

Hence,  $y=1$  is predicted for  $X$ .

Similarly for a point  $p$ , we get  $h(\theta) = -0.5 \leq 0$ . Hence  $y=0$  for this point.

In this e.g. we conclude that for points near  $l^{(1)}$  and  $l^{(2)}$ ,  $y=1$  and for points far away from  $l^{(1)}$  and  $l^{(2)}$ ,  $y=0$ .

In this way, we can learn pretty non linear boundaries.

Q. How do we chose these landmarks and what other similarity

functions can we use other than the Gaussian Kernel fn?

→ Given  $(n^{(1)}, y^{(1)}), (n^{(2)}, y^{(2)}), \dots (n^{(m)}, y^{(m)})$  {for m training examples}

choose  $l^{(1)} = n^{(2)}, l^{(2)} = n^{(3)}, \dots, l^{(m)} = n^{(m)}$

Given example  $n$ :

$f_1 = \text{similarity}(n, l^{(1)})$

$f_2 = \text{similarity}(n, l^{(2)})$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_m \end{bmatrix} \quad f_0 = 1$$

for training example  $(x^{(i)}, y^{(i)})$ :

$$f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)})$$

$$f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)})$$

$$\vdots \quad \leftarrow \quad f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) = \exp\left(-\frac{\|x^{(i)} - l^{(1)}\|^2}{2\sigma^2}\right) = 1$$

$$f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)})$$

Hence, one of these features = 1

Now  $x^{(i)}$  can be represented as  $f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$

Hypothesis: Given  $x_i$ , compute features  $f \in \mathbb{R}^{m+1}$

Predict " $y=1$ " if  $\theta^T f \geq 0$ .

$$\rightarrow = \theta_0 f_0 + \theta_1 f_1 + \theta_2 f_2 + \dots + \theta_m f_m \quad \{\theta \in \mathbb{R}^{m+1}\}$$

Training:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^m \theta_j^2$$

Solving this we  
(Solving  $\theta$  upto  $m$ ) get parameters for  
SVM.

$$\sum_j \theta_j^2 = \theta^T \theta, \quad \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \quad (\text{ignoring } \theta_0 \text{ as we do while regularize.})$$

instead of minimizing  $\|\theta\|^2$ , SVM minimizes  $\theta^T M \theta$   
 ↳ some matrix to take care of computational efficiencies, like if no. of training eg. is too large then modified helps in computation.

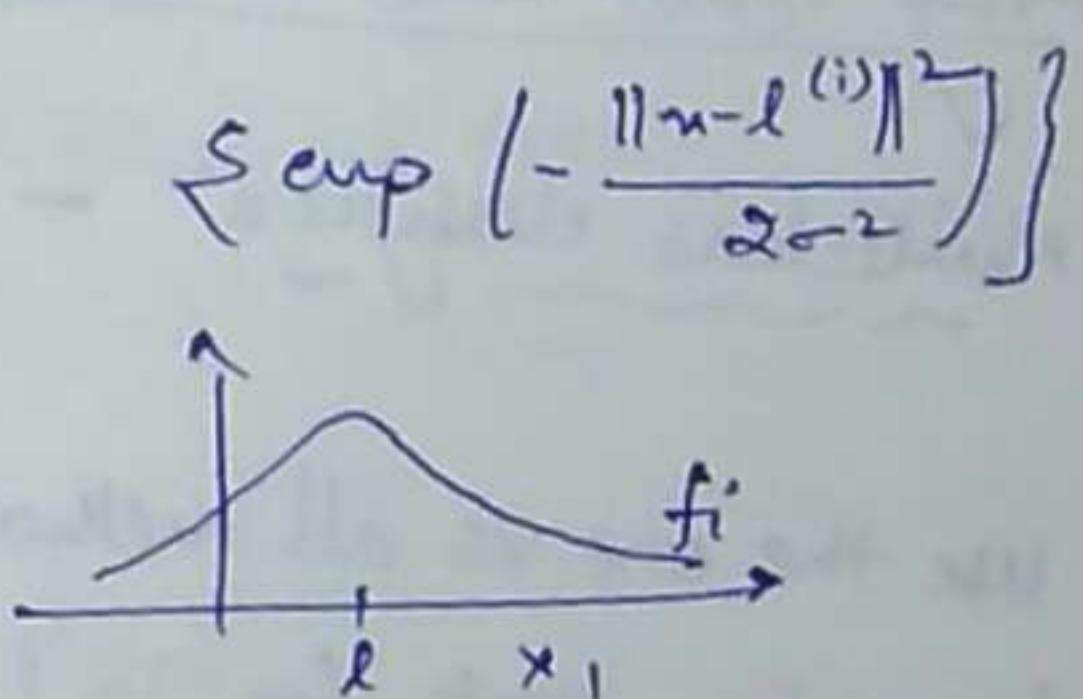
Note SVM and kernels go well together but logistic regression would run very slowly with kernels.

Q How to choose  $C$  for SVM? ,  $\sigma^2$  for SVM?

for optimization purpose.

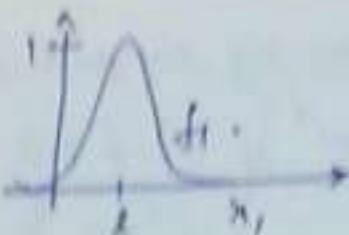
Ans large  $C$  : lower bias & high variance (small  $\lambda$ ) {prone to overfitting}  
small  $C$  : higher bias & lower variance. (large  $\lambda$ ) {prone to underfitting}

large  $\sigma^2$  : features  $f_i$  vary more smoothly  
 - higher bias, lower variance bcoz Gaussian Kernel that varies smoothly tends to give a hypothesis that varies slowly and smoothly



Gaussian Kernel tends to fall off slowly when  $\sigma^2$  is large.

Small  $\sigma^2$ : features vary less smoothly  
lower bias, higher variance.



e.g. No Kernel ("linear kernel")

Predict  $y=1$  if  $\theta^T x > 0$

This kernel is generally chosen when  $n$  is large and  $m$  is small to avoid overfitting.

e.g. Gaussian Kernel ( $f(x) = e^{-\frac{\|x - \mu\|^2}{2\sigma^2}}$ )  
where  $\mu^{(1)} = \mu^{(0)}$ .

This kernel is used when  $n$  is small but  $m$  is large to fit a more complex decision boundary.

Note To perform feature scaling before using the Gaussian Kernel, bcz otherwise some features with larger values would be given more after/weight than other features with smaller values. e.g. size of house would be given more weightage than no. of bedrooms if scaling is not done.

Not all similarity fn, similarity( $x, l$ ) make valid kernels. The kernels need to satisfy technical condition called "Mercer's theorem" to make sure that SVM packages optimize them correctly to calculate the parameters.

Other Kernels:

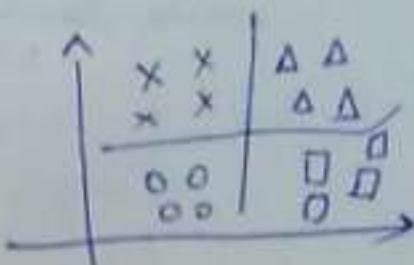
→ Polynomial kernel:  $K(x, l) = (x^T l)^2$  or  $(x^T l)^3$ ,  $(x^T l + 1)^3$ , etc.

The constant added and the degree of polynomial vary from 1 type of polynomial kernel to another.

→ String kernel, chi square kernel,  
histogram kernel, intersection kernel, ...

String kernel used when  $x$  is text string.

⇒ Multiclass classifier - SVM classifiers have generally built in multiclass classification functionality. Other method is to use the one vs all method where  $K$  SVMs are trained to distinguish  $y=i$  from the rest for  $i=1, 2, \dots, K$  to get  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$ . Then the class with largest  $(\theta^{(i)})^T x$  is picked up.



Ques When to use SVM and when to use logistic regression?

Ans  $n = \text{no. of features}$  ( $n \in \mathbb{R}^{n \times 1}$ ),  $m = \text{no. of training examples}$ .

- If  $n$  is large (relative to  $m$ ) : use LR or SVM without kernel ("linear")  
e.g. in text classification problem

beacuz when  $m$  is small, we don't have much data and so linear fn would probably do fine as we need not fit a very complicated non-linear function.

- If  $n$  is small,  $m$  is intermediate : Use SVM with Gaussian Kernel.  
( $n = 1-1000$ ,  $m = 10-10^3$ )

- If  $n$  is small,  $m$  is large : Create or add more features, then use LR or linear SVM, as in this case Gaussian SVM would be slow to run.

Note LR and linear SVM are quite similar. NN will work well for

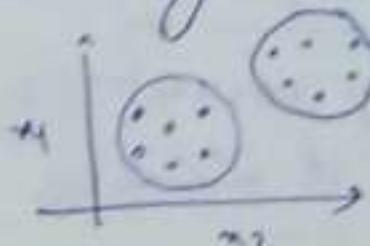
most of the above settings but they may be slower to train. The optimization problem that SVM has is a convex optimization problem and so good SVM optimizer software packages will always find a global minimum or something close to it, so for SVM we need not worry about local optima and for NN also, local optima is not a problem as it always finds global optima.

So, SVM : a very powerful tool to learn complex non-linear fn.

Linear SVM gives linear decision boundary just like LR does.

→ WEEK 8 :

→ Clustering : An Unsupervised learning algorithm.



Training set:  $\{n^{(1)}, n^{(2)}, n^{(3)}, \dots, n^{(m)}\}$

We need an algo to find some structure in the data for us.

An algo that finds clusters like the ones just circles is called clustering algorithm. Its applica<sup>n</sup> areas are:

- 1) Market segmentation
- 2) Social network analysis
- 3) Astronomical data analysis.

→ K-Means algorithm: a specific clustering algorithm. (most popular)

5/6/18 Skips:

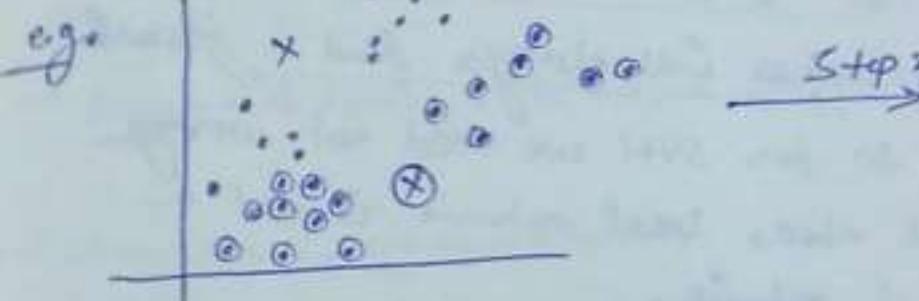
1) Randomly initialize two points called cluster centroids, when 3 want to group data into two clusters.

K-Means is an "iterative algorithm" and "it does two things -

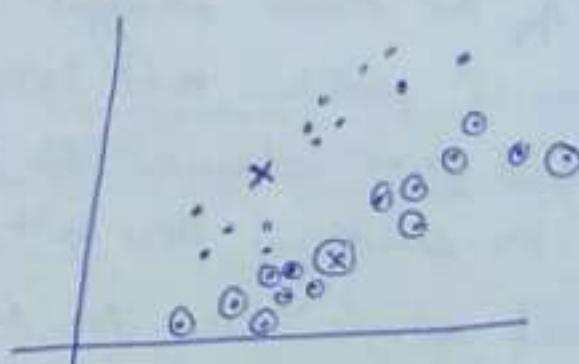
1) Cluster assignment step  
It checks all the datapoints and depending upon to which centroid the data point is closer, that cluster is assigned to it.

2) Move centroid step

The centroids are moved to the average of the points assigned to the same cluster.



① and x are cluster centroids randomly initialized. Data points close to ① are marked as ② and data points close to x are marked as .



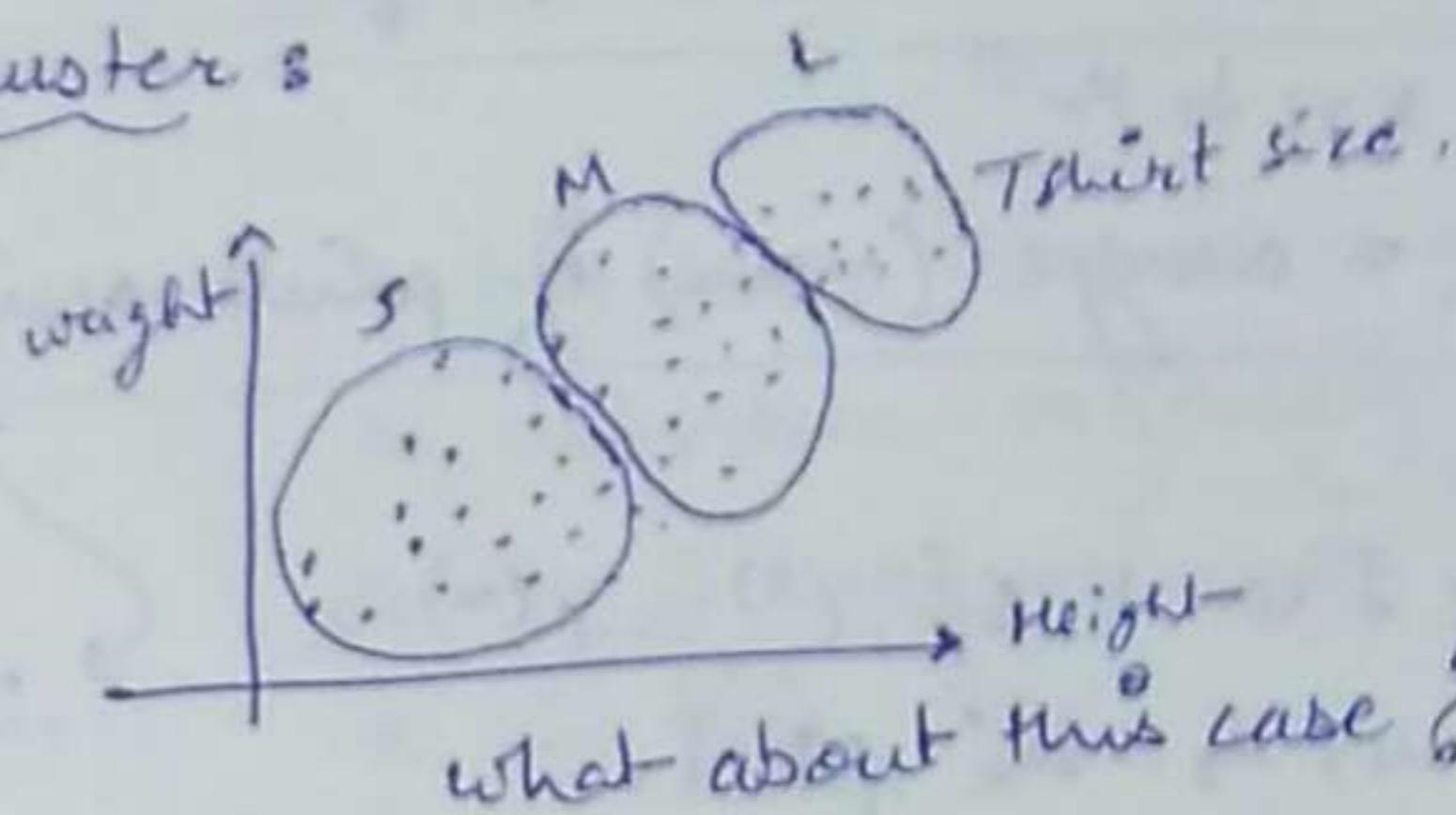
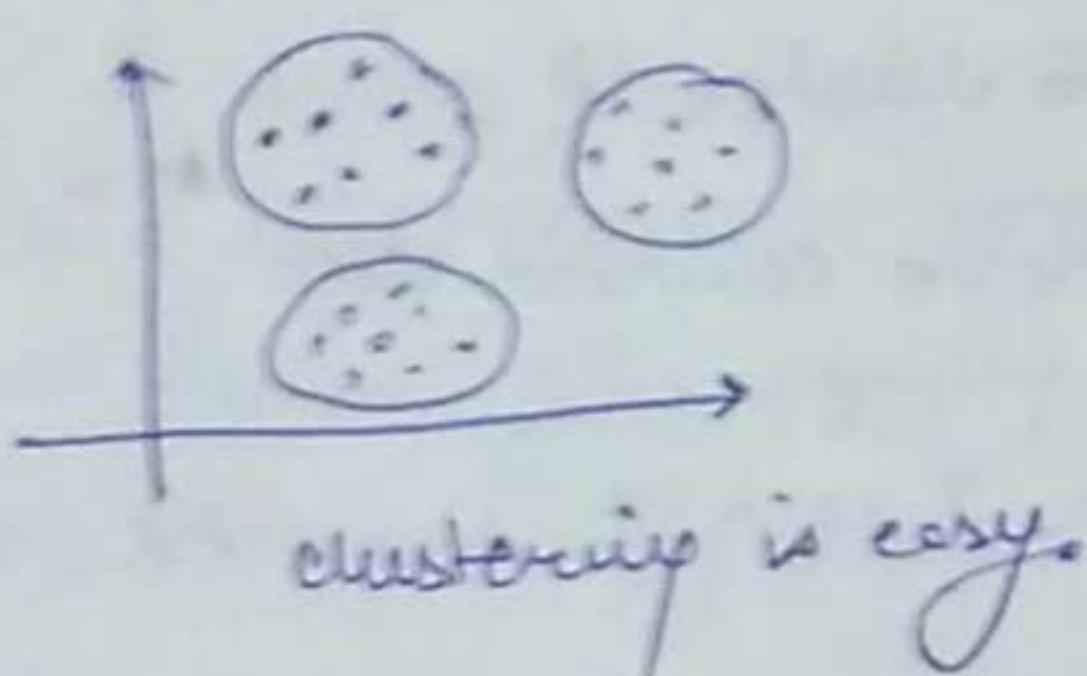
x ① and ② are moved to new pos<sup>n</sup> which are actually the avg of the loc<sup>n</sup> of . and ② points respectively.

Now Step 1 and 2 i.e. cluster assignment and centroid movement are repeated until the centroids stop changing their position further. At this point K-means converges giving two reasonably correct clusters. Each datapoint is assigned a cluster such that it minimizes

$\|x^{(i)} - \mu_k\|^2$  re distance b/w datapoint and cluster centroid.

If there arises a condn? where no datapoint is assigned to a particular cluster, then that particular cluster is eliminated or is reinitialized in case there is a need for all K clusters.

K-Means for non-separated clusters:



what about this case?

- suppose we want to make 3 clusters of Tshirt sizes 'S', 'M' and 'L'.  
(e.g. of market segmentation)

K-means optimize^ objective:

$c^{(i)}$  = index of cluster ( $1, 2, \dots, K$ ) to which example  $x^{(i)}$  is currently assigned.

$\mu_k$  = cluster centroid  $k$  ( $\mu_k \in R^n$ ) {location}

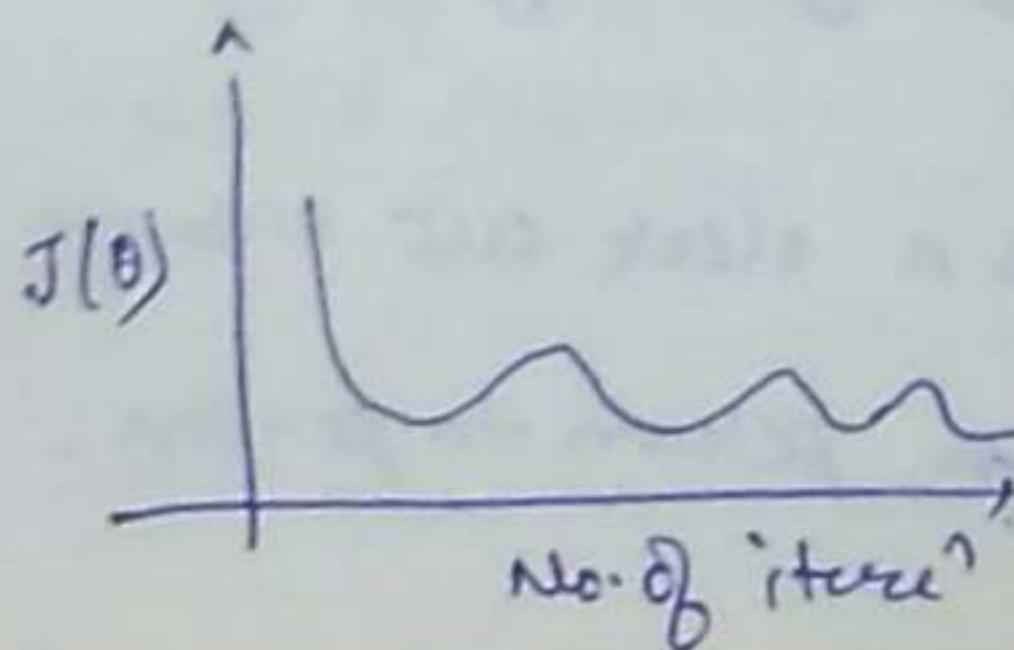
$K$  = total no. of clusters.

$\mu_{c^{(i)}}$  : cluster centroid of cluster to which example  $x^{(i)}$  has been assigned.

Optimize^ objective :  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$

$$\min_{\substack{c^{(1)}, \dots, c^{(m)}, \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

This cost  $J$  is also called distortion cost  $J$  or the distor^ of K-Means alg.



Such a  $J(\theta)$  is not possible w.r.t. iterations. There must be some bug in the code if such  $J(\theta)$  is obtained.

## K-means algorithm:

Randomly initialize K cluster centroids  $u_1, u_2, \dots, u_K \in \mathbb{R}^n$

Repeat {

for  $i = 1$  to  $m$

$c^{(i)} = \text{index}(\text{from}, \text{to } K)$  of cluster centroid closest to  $x^{(i)}$

for  $k = 1$  to  $K$

$u_k = \text{average (mean)} \text{ of points assigned to cluster } k.$

Cluster assignment step

Move centroid step.

minimize  $J$  w.r.t  $c^{(1)}, c^{(2)}, \dots, c^{(m)}$   
holding  $u_1, u_2, \dots, u_K$  fixed.

minimize  $J$  w.r.t.  $u_1, \dots, u_K$

Q. How to initialize K means and avoid local optima as well?

Ans Should have  $K < m$

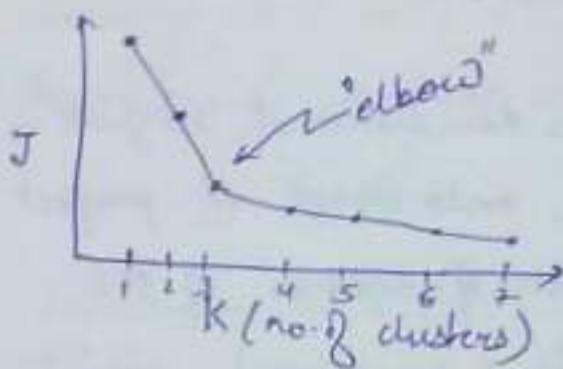
Randomly pick K training examples  
Set  $u_1, \dots, u_K$  equal to these K examples.] i.e.  $u_1 = x^{(i)}$  if  $K=2$   
 $u_2 = x^{(j)}$  ; ; are some random training exs

Depending on the random initialize, K means can end up with different solutions. So, to ensure that K means end up with global and not local optima, try multiple random initializations. Then calculate the cost corresponding to each initialize and pick the one with minimum cost. This works even better when K is small (i.e 2 to 10). If K is very large, the 1st random initialize itself would lead to pretty good result and many random initialize won't make much difference.

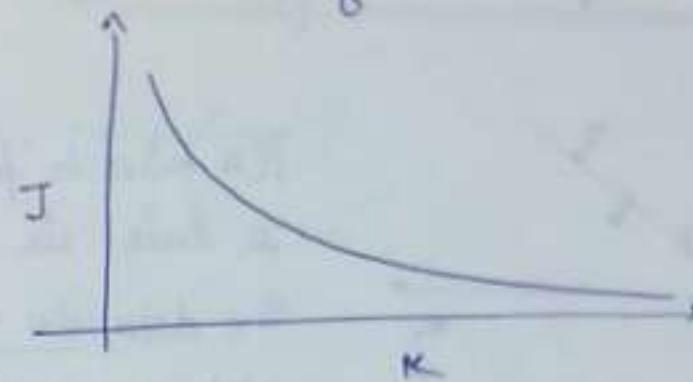
Q How to choose the no. of clusters K?

Ans The most common way is to choose it manually by looking at visualizations and K-means clustering results.  
Being unsupervised learning, there isn't always a clear cut answer to the possible no. of clusters as it depends from person to person.

- elbow method: to find a reasonable value of  $K$ .



It seems like  $K=3$  would be a right decision to make, as distortion goes rapidly down till 3 but hardly changes after 3.



In this case, there is no clear elbow as distortion continually goes down. So, it is hard to make a choice for  $K$  in this case.

Note If  $J$  is much higher for  $k=5$  than for  $k=3$ , that means in the run with  $K=5$ , k-means got stuck in a bad local minimum and k-means should be rerun with multiple random initializations.

- Sometimes we run K-means to get cluster to use for some later purpose - e.g. in market segmentation (T-shirt e.g. taken earlier). In this case, evaluate K-Means based on a metric for how well it performs for that later purpose.

\* **Dimensionality Reduction**: Another unsupervised learning algorithm needed for

Data compression  
(helps speed learning algo)  
(helps save memory)

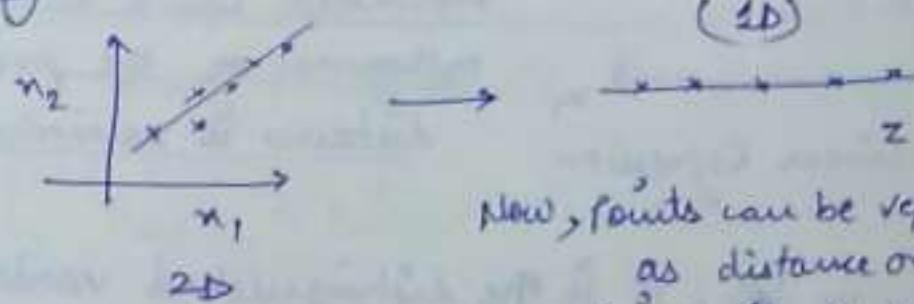
Visualization  
(reduce data to 2D or 3D in order to visualize it.)  
( $K=2$  or 3 while visualization.)

as we can maximum plot 3 features.

( $n_1, n_2$  are features)

This dimensionality reduction is done by removing redundant or less important features or combining features to give one feature.

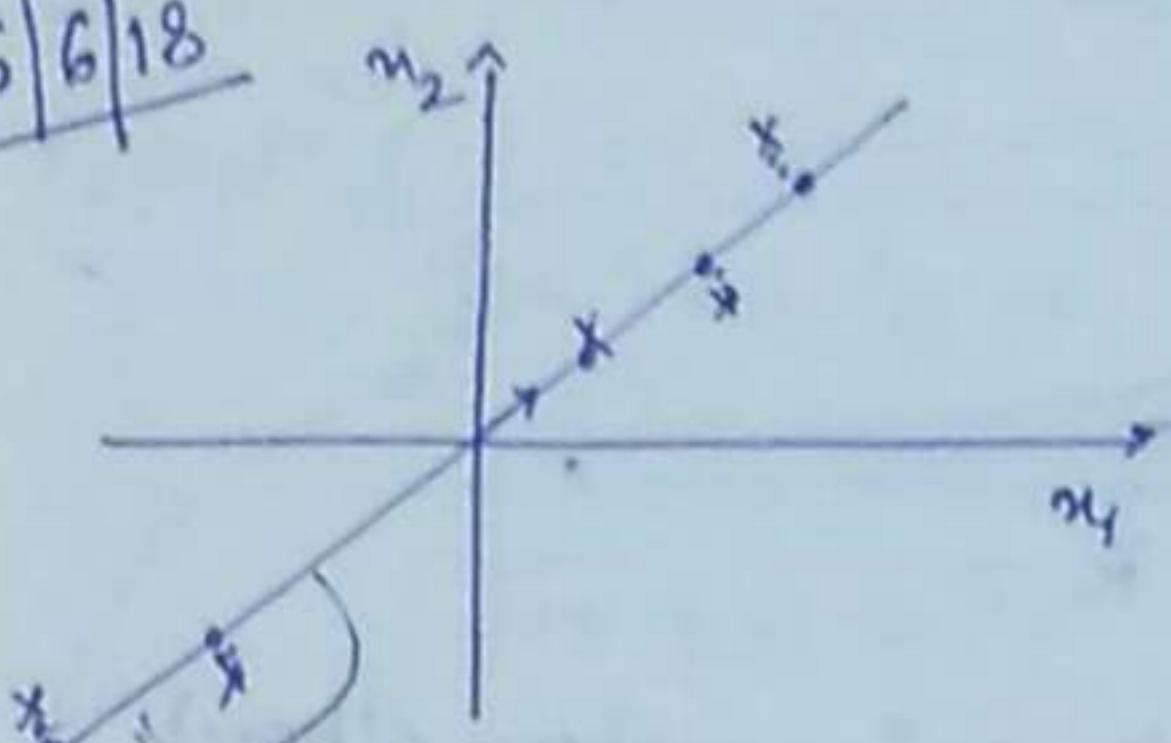
e.g. 2D to 1D.



Now, points can be represented as distance on this line. So,  $n_1$  and  $n_2$  are reduced to just 1 number  $z$ .

→ Principal Component Analysis (PCA) : a dimensionality reduction algorithm.

6/6/18



x: actual data points

•: projected data points

→ best line that can minimize the projection error.

PCA tries to find a lower dimensional surface, a line in this case, onto which to project the data so that the sum of squares of the little dotted line segments is minimized, i.e. it tries to minimize the projection error.

Before applying, it is a standard practice to perform mean normalization & feature scaling so that features  $n_1$  and  $n_2$  have mean 0 and should have comparable ranges of values.

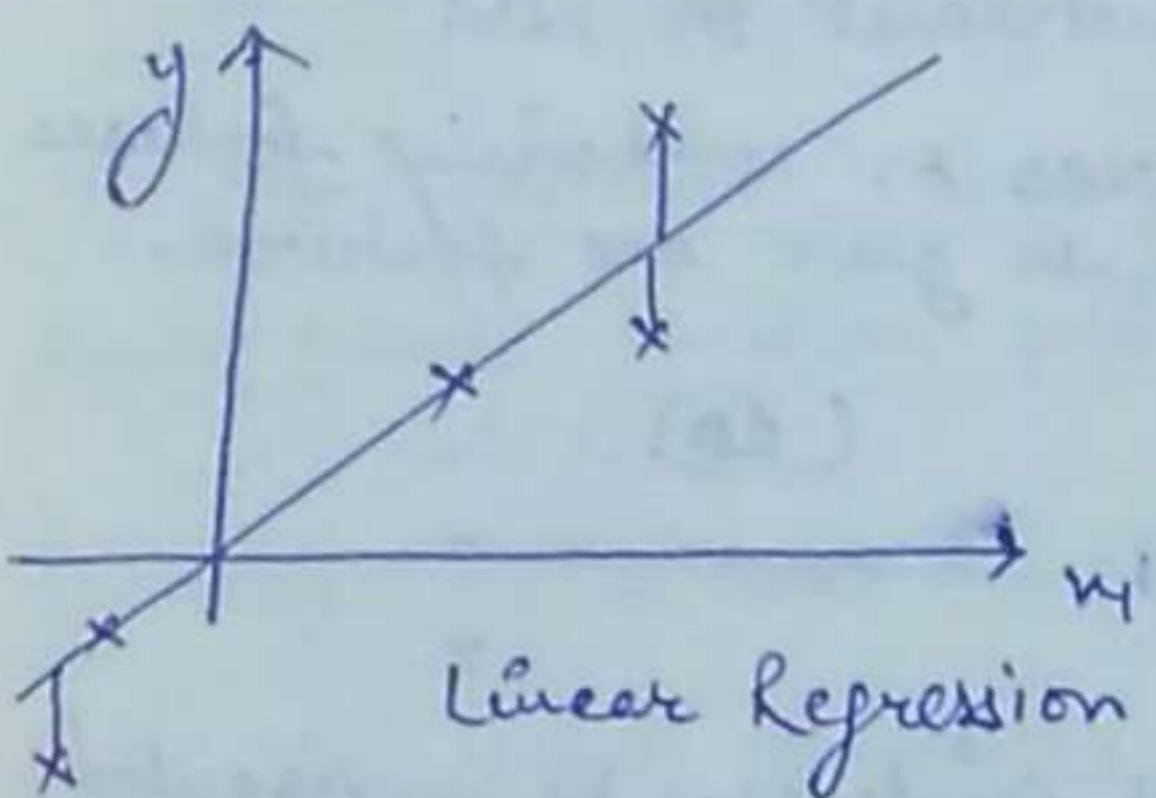
So, PCA finds a direction (a vector  $u^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize projection error.

Generally, we need to reduce  $n$ -dimensional to  $k$ -dimensional data where  $k$  vectors are found on which to project the data.

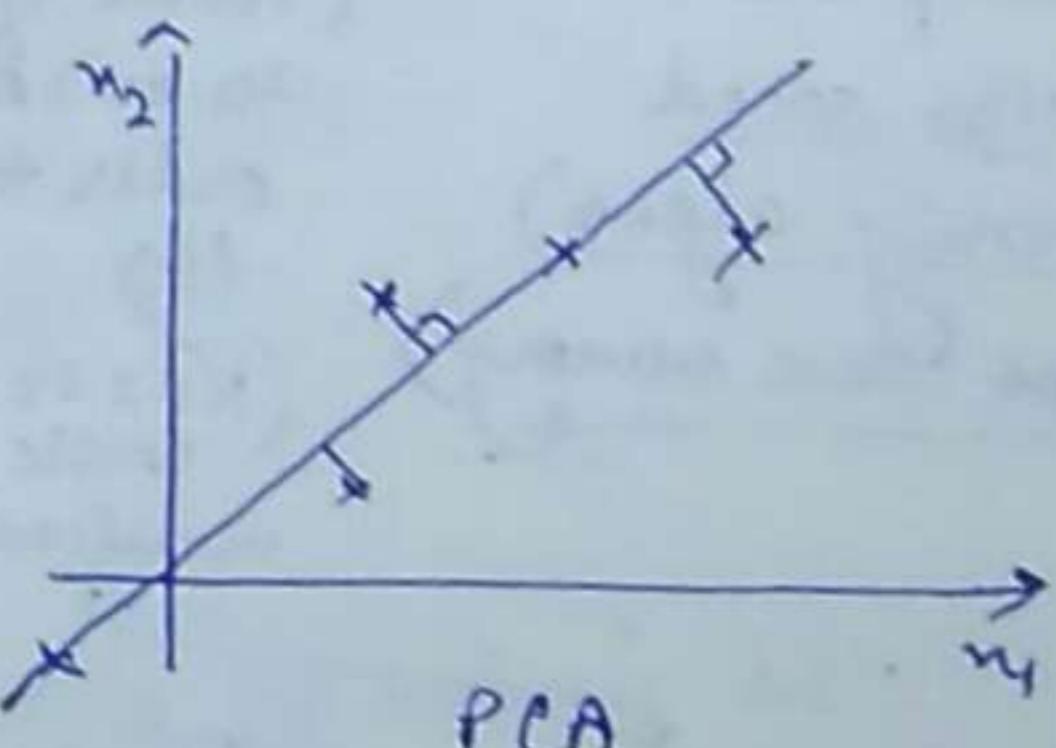
Q How PCA relates to linear regression?

Ans PCA and linear regression are 2 completely different algo though they look the same.

In linear regression, we try fit a straight line so as to minimize the square error b/w point and this straight line.



In linear regression, the vertical distance is being minimized while in PCA orthogonal or the project distance is minimized.



Moreover, in LR, it is the distinguished variable 'y' we are trying to predict whereas in PCA there is no such  $y$  we are trying to predict rather we have a list of features  $n_1, n_2 \dots n_r$  which are all treated equally.

To apply PCA, we 1<sup>st</sup> need to do data preprocessing :

Training set :  $x^{(1)}, x^{(2)}, \dots x^{(m)}$

Preprocessing (feature scaling / mean normalize)

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each  $x_j^{(i)}$  with  $x_j - \mu_j$

of different features on different scales (e.g.  $x_1$  = size of house,  $x_2$  = # bedrooms)  
scale features to have comparable range of values.

$$x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{s_j}$$

2D to 1D reduce by PCA means:  $x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}^1$

3D to 2D . . . . . :  $x^{(i)} \in \mathbb{R}^3 \rightarrow z^{(i)} \in \mathbb{R}^2$  and so on.  
where  $z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$  for 2D.

So, PCA needs to calculate  $u$  (i.e vectors on which to project) and  $z$  (i.e new data points in reduced dimensions)

PCA also to reduce a n-D to k-D:

Compute covariance matrix (n x n matrix)

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \quad \{ \Sigma \text{ is Greek symbol } \Sigma \text{ to denote matrix} \}$$

Compute eigen vectors of matrix  $\Sigma$  -

$$[U, S, V] = \text{svd}(\Sigma) \quad \{ \text{Octave command for singular value decomposition} \}$$

Note covariance matrix always satisfies a property 'symmetric +ve definite' and  
so  $\text{svd}(\Sigma) = \text{eig}(\Sigma)$ . Otherwise eig & svd are two very different functions.

$$U = \left[ \begin{array}{c|c|c|c} 1 & & & \\ u^{(1)} & u^{(2)} & \cdots & u^{(m)} \\ \hline 1 & & & \end{array} \right] \quad U \in \mathbb{R}^{n \times n}$$

and to reduce to k-dimensions, we take the first k columns i.e  $u^{(1)} \text{ to } u^{(k)}$ .

$$U_{\text{reduce}} = \left[ \begin{array}{c|c|c|c} 1 & & & \\ u^{(1)} & u^{(2)} & \cdots & u^{(k)} \\ \hline 1 & & & \end{array} \right] \quad U_{\text{reduce}} \in \mathbb{R}^{n \times k}$$

$$z = U_{\text{reduce}}^T x \in \mathbb{R}^{k \times 1} \quad \{ \text{Hence, we got the k vectors to project our data on} \}$$

$$z^{(i)} = U_{\text{reduce}}^T x^{(i)} \quad (\text{for particular training e.g.})$$

the mathematical proof that PCA gives optimised reduced dimensions is quite complex.

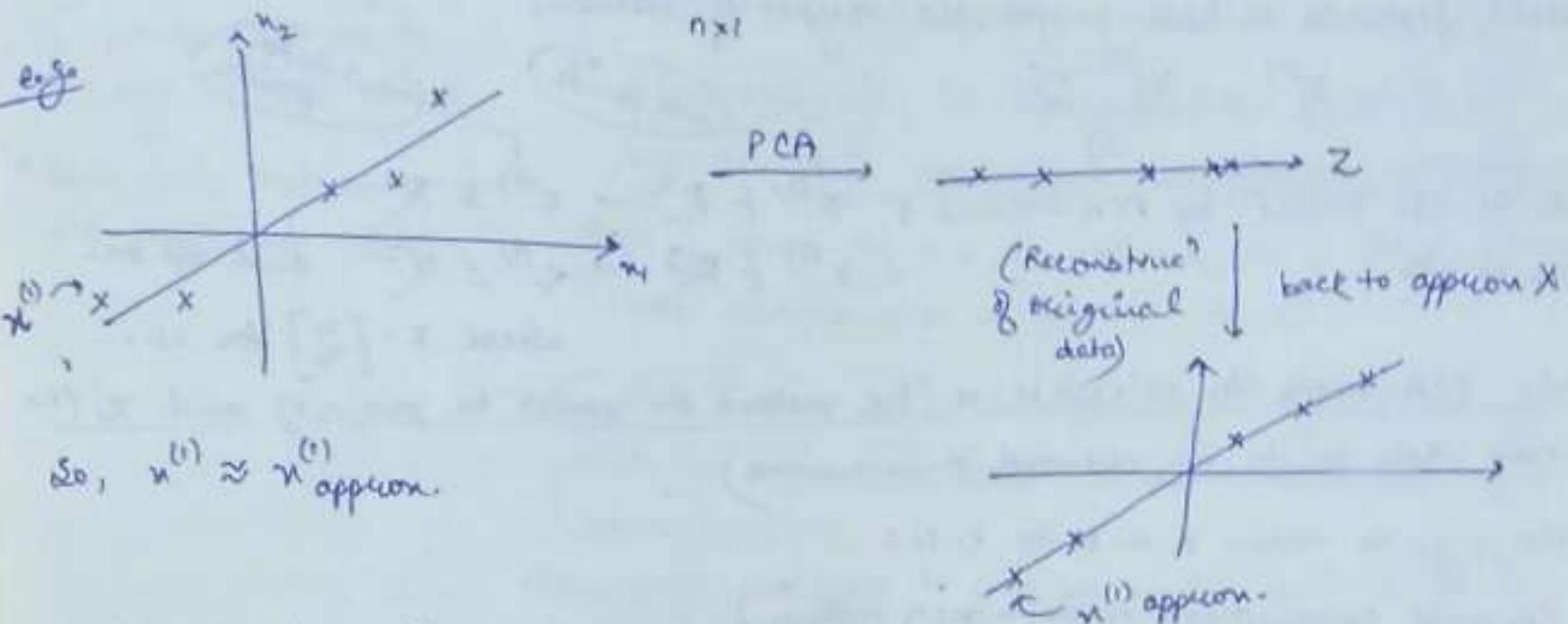
Q How to go back from  $z_i$  (compressed representation) to  $x_i$  (original representation)

Ans We know  $z = U_{\text{reduce}}^T X$

and so,  $X_{\text{approx}} = \underbrace{U_{\text{reduce}}}_{n \times K} \underbrace{\cdot z}_{K \times 1}$ ,  
 $n \times 1$

$X_{\text{approx}}$  is closer to original  $X$  from where we started.

e.g.



Q How to choose  $K$ , i.e. the dimension to which  $X$  need to be reduced?

Ans  $K \rightarrow$  parameter of PCA algorithm.  
or no. of principal components that we've retained.

Avg. squared projection error =  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

Total variance in the data :  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

(on avg. how far are our training examples from vectors, from just being all zeros)

Typically choose  $k$  to be the smallest value so that:

$$\textcircled{1} \quad \frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} < 0.01$$

i.e (1%)

\*\*

✓ "99% of variance is retained" {a good way to repeat the k chosen?}

In real life applica<sup>n</sup>, many of the features are highly correlated and so it's quick easy to reduce the dimensions significantly while retaining 99% variance.

$S \in [U, S, V]$  from  $\text{svd}(\text{sigma})$  is a  $n \times n$  matrix which is diagonal.

$$S = \begin{bmatrix} s_{11} & & 0 \\ & s_{22} & \\ 0 & & \ddots s_{nn} \end{bmatrix}$$

for given  $K$ ,  $\Theta$  can be calculated as:

$$\Theta = \begin{bmatrix} 1 - \frac{\sum_{i=1}^K s_{ii}}{\sum_{i=1}^n s_{ii}} \end{bmatrix}$$

(much simpler expression than ①)

We need to get just once the value of  $S$  to find the optimum  $K$ . We need not run PCA again & again from scratch as  $\Theta$  has been obtained from  $\Theta$ .

so, we need to check for which smallest value of  $K$ ,  $\Theta \leq 0.01$

Q How PCA helps to speed up a learning algorithm?

Ans ① Supervised learning speedup

let's say we have this data:  $(u^{(1)}, y^{(1)}), (u^{(2)}, y^{(2)}), \dots, (u^{(m)}, y^{(m)})$  and say  $u^{(i)} \in \mathbb{R}^{10K}$  i.e. we have very high dimensional data. such huge data will surely slow down any algo be it LR or SVM or some other. Hence, we need PCA so as to speed it up with lesser features.

Procedure:

Extract inputs: i.e unlabeled dataset  $\underbrace{(u^{(1)}, u^{(2)}, \dots, u^{(m)})}_{\downarrow \text{PCA}} \in \mathbb{R}^{10K}$

$$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{1000}$$

so, New training set:  $(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}) \dots, (z^{(m)}, y^{(m)})$

→ lower dimensional feature vector than  $u$ .  
(obtained through PCA)

Now, this reduced dimensional data can be fed into our learning algorithm (SVM, LR, NN, etc) to have  $h_0(z) = \frac{1}{1+e^{-\theta^T z}}$ . Any new example can be mapped using PCA to give reduced feature vector  $z$ .

Note) Mapping  $x^{(i)} \rightarrow z^{(i)}$  should be defined by running PCA only on the training set. This mapping can be applied as well to the e.g.  $x_{cv}^{(i)}$  and  $x_{test}^{(i)}$ . This is so becuz, while mapping Unseen is calculated mean normalized and feature scaling is done which are actually parameter learned by PCA, and so this can be done on trainset only.

→ Bad use of PCA: to prevent overfitting

PCA might work OK but use regularization instead.

becuz PCA does not use labels  $y$ , it only sees  $x$  and reduces it to lower dimension, hence with probability of throwing away valuable information which regularization won't do as it takes into account the labels  $y$  unlike PCA.

\* Before implementing PLA, first try running with the original/raw data  $x^{(i)}$ . Only if that doesn't work, implement PCA and consider using  $z^{(i)}$ . when you think there is a speed or memory problem, think of applying PCA, otherwise try doing with the original data only.

\* WEEK 9 :

Anomaly detection - (Supervised & Unsupervised learning algo)

$x^{(i)}$  = features of user i's activities

Model  $p(x)$  from data

Identify unusual users by checking which have  $p(x) < \epsilon$

If there are too many false pos detected, ↓  $\epsilon$ .

anomalous

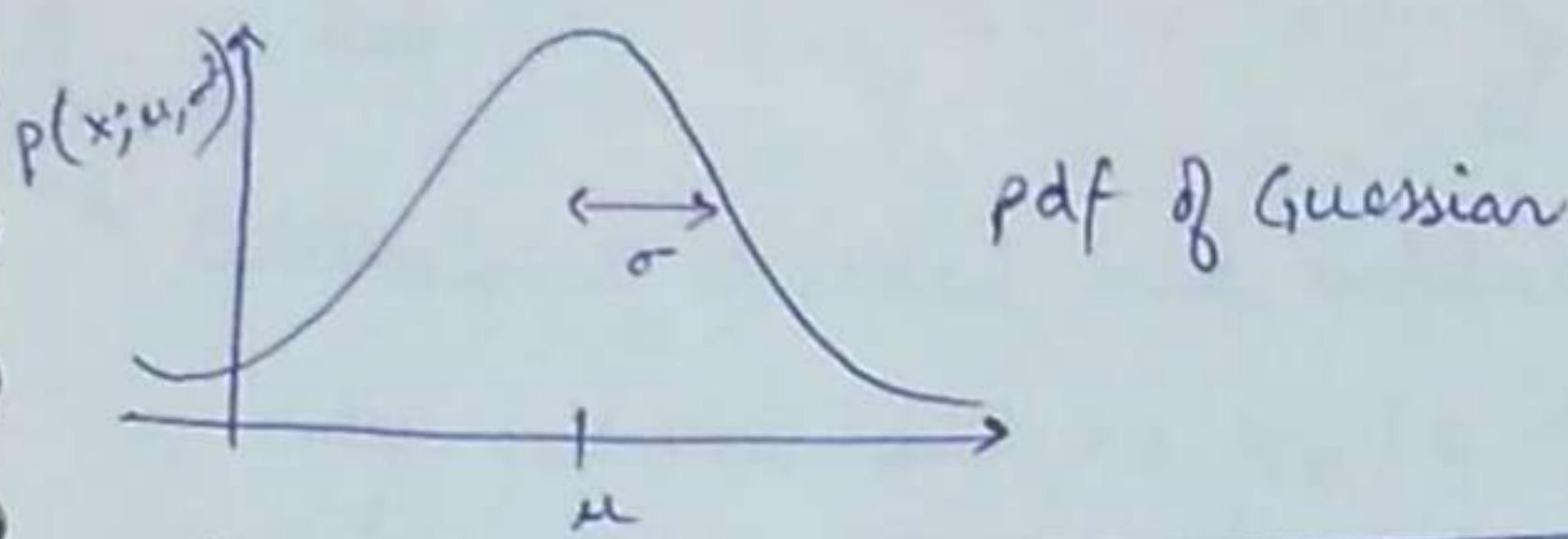
} applies in fraud detection, manufacturing or monitoring computers in a data center, etc.

Gaussian Distribution (Normal distribution) : to develop an anomaly detect algo.

Say  $n \sim \mathcal{N}$ , if  $n$  is a distributed Gaussian with mean  $\mu$ , variance  $\sigma^2$ .

$X \sim \mathcal{N}(\mu, \sigma^2)$   
 "distributed as"  $\rightarrow$  Normal distribu<sup>n</sup>  
 (same as Gaussian)

$\mu$ : mean  
 $\sigma^2$ : variance  
 $\sigma$ : standard deviation



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

Probability of  $n$  parameterized by  $2$  parameters  $\mu$  &  $\sigma^2$ .

### - Parameter estimation:

Dataset:  $\{n^{(1)}, n^{(2)}, \dots, n^{(m)}\}$   $n^{(i)} \in \mathbb{R}$

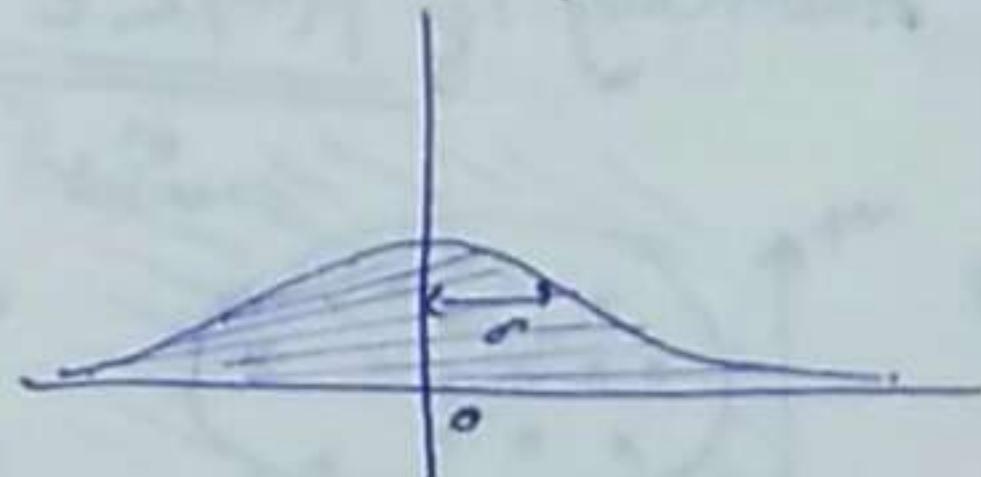
Now, we need to find  $\mu, \sigma$  such that this dataset fits  $\mathcal{N}(\mu, \sigma^2)$ .

So,  $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$  (means center of the dataset distribu<sup>n</sup>)

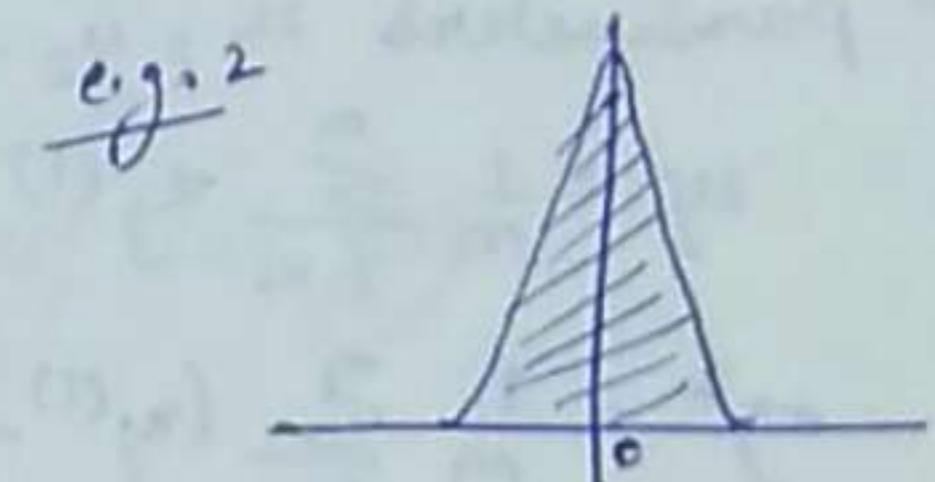
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

Here,  $\mu$  and  $\sigma^2$  are actually the mean likelihood estimates of the parameters  $\mu$  and  $\sigma^2$ .

The area in both the examples integrate to 1 as it's the property of probability distribu<sup>n</sup>.

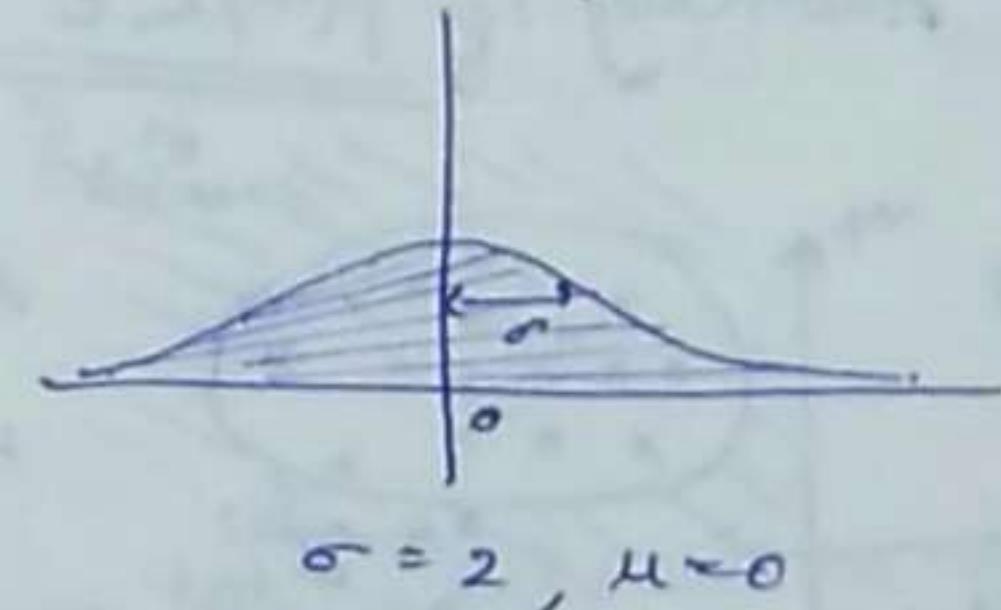


$$\mu=0, \sigma=1$$



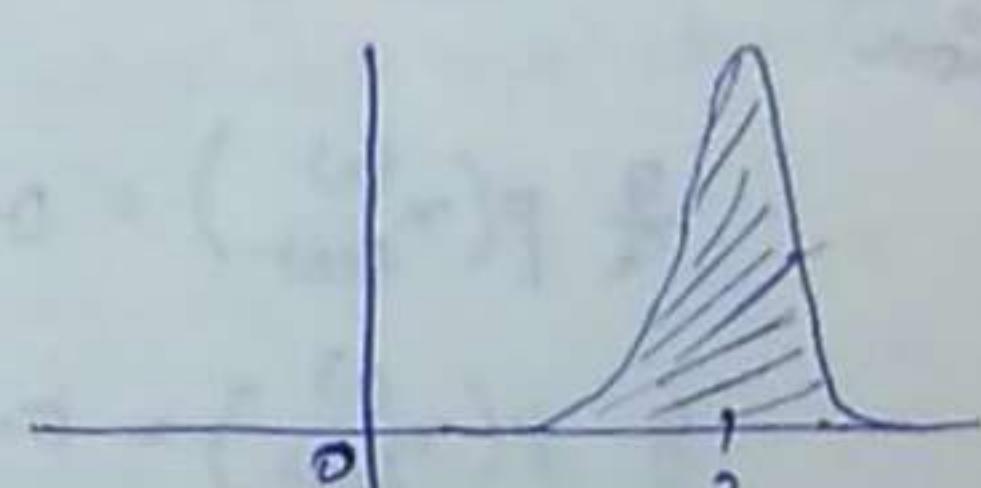
$$\mu=0, \sigma=0.5$$

e.g. 3



$$\sigma=2, \mu=0$$

e.g. 4



$$\mu=3, \sigma=0.5$$

### 7/6/18 Density Estimation problem:

Training set:  $\{n^{(1)}, \dots, n^{(m)}\}$

Each example  $n$  is  $\mathbb{R}^n$

We will now model  $p(n)$  from dataset to figure out imp & less imp features.

$$p(n) = p(n_1; \mu_1, \sigma_1^2) p(n_2; \mu_2, \sigma_2^2) p(n_3; \mu_3, \sigma_3^2) \dots p(n_n; \mu_n, \sigma_n^2)$$

We assume that  $n_i$  is distributed according to some Gaussian distribution as  $n_i \sim \mathcal{N}(\mu, \sigma^2)$  and so on for all features.

(i) The eq' p( $\mathbf{x}$ ) corresponds to independence assumption on the values of the features  $x_1$  through  $x_n$ , but the model works for even when the features are not close to independent.

$$p(\mathbf{x}) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

- Anomaly detection algo:

- 1) Choose features  $x_i$  that you think might be indicative of anomalous eg.
- 2) Fit parameters  $\mu_1, \mu_2, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

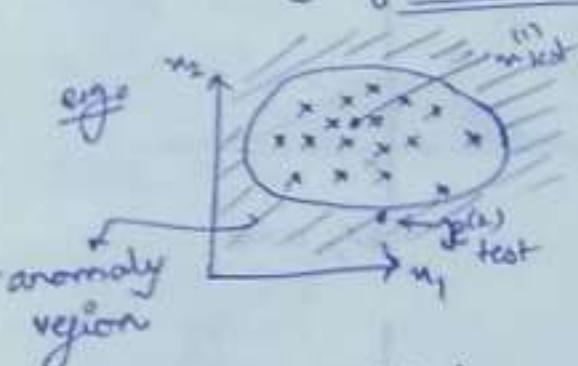
$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \rightarrow \mu = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \quad (\text{vectorized version})$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

- 3) Given new example  $\mathbf{x}$ , compute  $p(\mathbf{x})$ ;

$$p(\mathbf{x}) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly: if  $p(\mathbf{x}) < \epsilon$ .



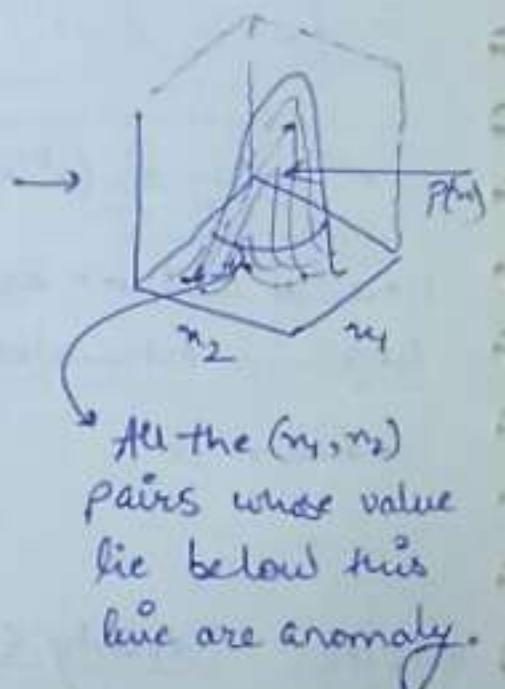
$\mu_1 = 5, \sigma_1 = 2$  (corresponding to  $n_1$ )

$\mu_2 = 3, \sigma_2 = 1$  ( \* " " to  $n_2$ )

Say  $\epsilon = 0.02$  (chosen value)

if  $p(x_{1 \text{ test}}^{(1)}) = 0.0015 < \epsilon$  Hence anomaly

if  $p(x_{2 \text{ test}}^{(1)}) = 0.0426 > \epsilon$  Hence, no anomaly



- Evaluation of anomaly detection algo:

Real no. evaluation is important while developing a learning algorithm (choosing features, etc) for easier/simpler decision making.

e.g. Aircraft engines:

10000 good engines } dataset (it's OK if some bad engines slip into the 20 flawed \* } 10K dataset, majority of them should be good.)

Training set: 6000 good entries. ( $y=0$ )  $\rightarrow$  fit these to  $p(n)$   
 CV: 2000 " " " ( $y=0$ ), 10 anomalous ( $y=1$ )  
 Test: " " " ( $y=0$ ), 10 anomalous ( $y=1$ )

Splitting of labelled and unlabelled e.g.s, with training set consisting of  
 only features and no  $y$ , i.e. training is set of unlabeled e.g.s.

Note The e.g. in cv and test set should be different.

After modelling  $p(n)$  on training set,

$\text{predict } y = \begin{cases} 1 & \text{if } p(n) < t \text{ (anomaly)} \\ 0 & \text{if } p(n) > t \text{ (normal)} \end{cases}$  on cv/test e.g.

Possible evaluation metrics: F<sub>1</sub>-score, Precision/Recall. by comparing the  
 predicted  $y$  with actual  $y$ .

$t$  can be chosen by running on cv.

Note Classification accuracy is not a good way to measure the accuracy the data  
 is skewed and any algo predicting more no. of  $y=0$  will come out as  
 a better algo.

Q Anomaly detection v/s supervised learning?

Ans In above eg we saw that we use labeled data (anomalous or not),  
 then why don't we simply use LR or NN algo. Why do we switch to  
 anomaly detection.

### Anomaly detection

- Very small no. of +ve examples (0-20), large no. of -ve examples ( $y=0$ ) bcoz to fit  $p(n)$  we need only -ve examples (aircrafts in good cond?).
- Many diff types of anomalies. Hard for any algorithm to learn from +ve e.g.s what the anomalies look like, future anomalies may look like any of the anomalous e.g.s we've seen so far.

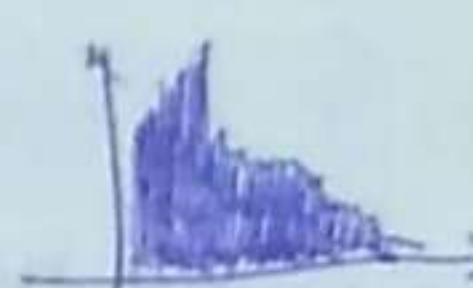
### Supervised learning

- Large no. of +ve and -ve e.g.s.
- Enough +ve e.g.s for algorithm to get a sense of what +ve e.g.s are like, future +ve e.g.s likely to be similar to ones in training set.

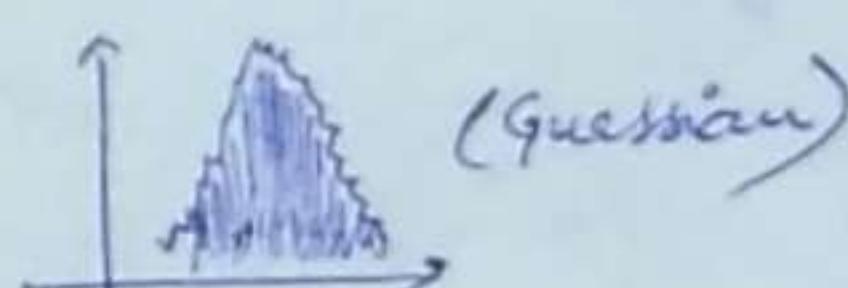
Spam email detec<sup>n</sup> can be done by supervised learning bcoz we have enough +ve examples available to capture all kind of spam emails and make the model learn which future mails would be spam and which won't be. Manufacturing defect problems are solved using anomaly detection bcoz in this case we don't have much +ve examples. Similarly, fraud detec<sup>n</sup> is anomaly detection problem as we don't have much data of fraud access to our system. If we do have, problem shifts to being a supervised learning problem.

Q. what features to choose to use in anomaly detection algorithm?

Ans we choose  $\sigma$ , i.e. considering data to be Gaussian distributed. Even if data isn't Gaussian, this works fine but still it's good to do a sanity check. If the plot doesn't look Gaussian,

e.g.  , what we do is play with different transformations of the data in order to make it more Gaussian.

$\downarrow \log(n)$  this would improve our results.



Any transformation can be done with any of the features to transform the data to Gaussian. e.g.  $n_1 \leftarrow \log(n)$

$$n_2 \leftarrow \log(n_1 + 1)$$

$$n_3 \leftarrow \sqrt{n_2} \text{, etc.}$$

Feature analysis can help find features for anomaly detection algorithms.

- We want  $p(x)$  to be large for normal examples  $x$ .

$p(x)$  to be small " anomalous examples  $x$ .

Most common problem:  $p(x)$  is comparable for both normal & anomalous e.g., thus failing to detect anomaly.

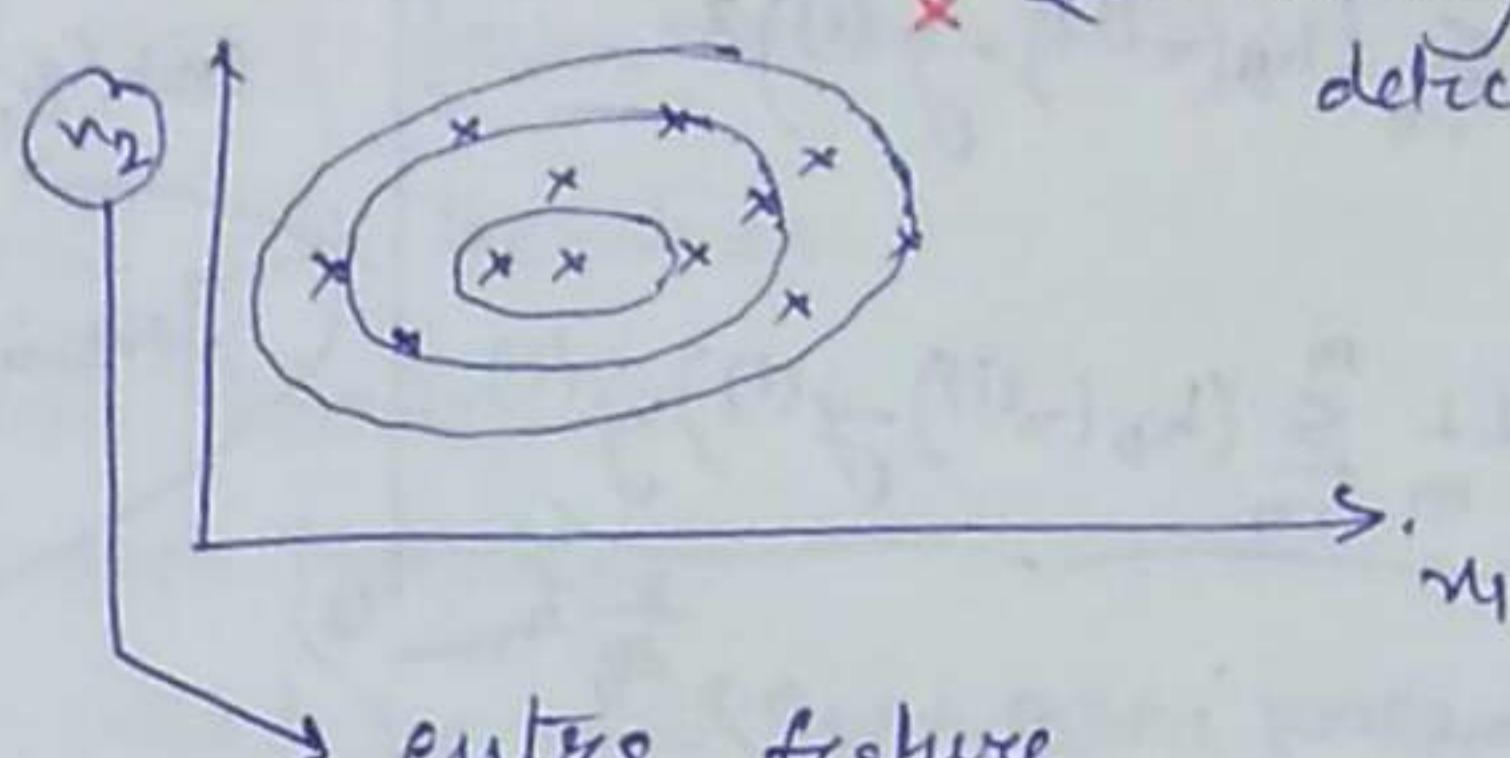
Soln: look at the anomaly that the algo is failing to flag and see if that inspires to create some new feature. So, find something unusual about that anomaly and use that to create a new feature so that with this new feature, it becomes easier to distinguish

the anomalies from the good examples.

e.g,



case failed to be detected (as it was in mid of the crowd of good e.g.s.)



entire feature added by studying the anomaly behaviour.

Note Choose features that might take on unusually large or small values in the event of an anomaly.

e.g. if  $n_1 = \text{weight}$ ,  $n_3 = \text{BMI}$  can be created to capture unusual health anomalies.

⇒ Multivariate Gaussian Distribution: to catch anomalies that the earlier algo didn't.

8/6/18

Note Best way to get accuracy is to take a low bias learning algo and train that on a lot of data.

large datasets pose large computational problems.

Two main ideas for viewing with very big datasets ↗ Stochastic Gradient Descent  
Map Reduce

⇒ Stochastic Gradient Descent: (SGD)

Gradient descent is computationally quite expensive with large datasets & so stochastic GD comes into picture. This is in replacement for the general batch gradient descent (where we're looking at all of the training examples at a time,) in case of large ' $m$ '. This SGD need not look at all training examples in one iteration rather it needs to look at only a single training example in one iteration.



### Batch GD

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

Repeat  $\mathcal{E}$

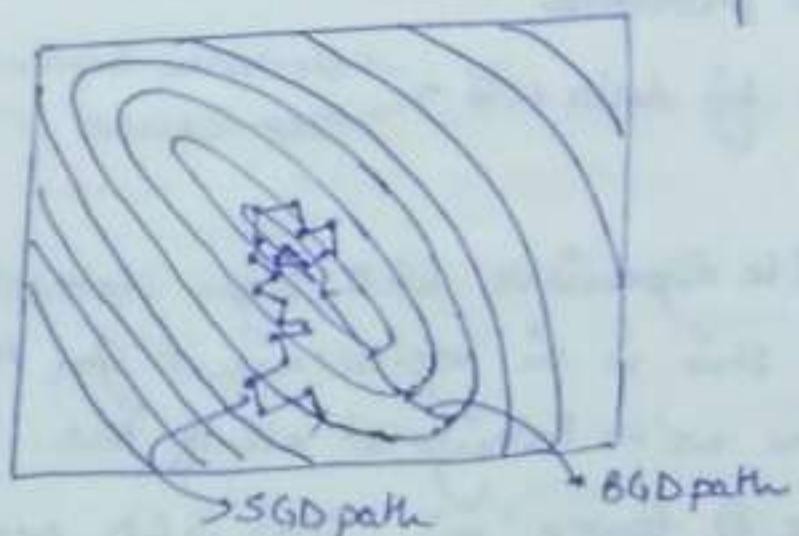
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(for every  $j = 0, \dots, n$ )

}

Here what is done is  $\theta_j$  is adjusted for all the training examples, then  $\theta_2$  is adjusted, then  $\theta_3$  and so on upto  $\theta_n$ .

Hence, in BGD to move towards the global minima we need to go through all the training exgs but in SGD each training example moves us a bit closer to the global minima.



Note  $J_{\text{train}}(\theta)$  should go down with every iteration of batch gradient descent (assuming a well tuned learning rate  $\alpha$ ) but not necessarily in SGD. SGD is but much faster than BGD for large  $m$ .

### SGD

$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

→ cost of parameter  $\theta$  w.r.t the training eg.  $(x^{(i)}, y^{(i)})$

Hence, this algo checks how much the hypothesis costs for training eg.  $(x^{(i)}, y^{(i)})$

Steps of SGD:

- i) Randomly shuffle dataset
- ii) Repeat  $\mathcal{E}$

for  $i = 1, \dots, m \in \mathcal{E}$

$$\theta_j := \theta_j - \alpha \frac{(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}{\partial \theta_j}$$

}

$$( \text{for } j = 0, \dots, n ) \quad \left. \frac{\partial \text{cost}(\theta, (x^{(i)}, y^{(i)}))}{\partial \theta_j} \right)$$

In this algo, we are basically trying to fit the parameters  $\theta$  for our 1<sup>st</sup> training eg, then 2<sup>nd</sup>, then 3<sup>rd</sup> and so on adjusting  $\theta$  until the  $m^{\text{th}}$  training eg.

⇒ Mini Batch Gradient Descent : (can sometimes work even a bit faster than SGD)

• Su BGD : use all  $m$  examples in each iteration.

• SGD : use 1 example in each iteration.

mini BGD : use ' $b$ ' examples in each iteration.

↳ parameter called the mini batch size. Typical range of  $b$  is 2 to 100.

Say  $b = 10$ ,  $m = 1000$ .

Repeat {

for  $i = 1, 11, 21, 31, \dots, 991$  {

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=1}^{10} (\text{log}(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

(for every  $j = 0, \dots, n$ )

}

This is faster than BGD as updates are now made after seeing 10 training egs rather than to wait for the entire training set to see an update.

Q: Why do we need to look at  $b$  examples at a time than to look at 1 training example as in SGD?

A: Due to 'Vectorization'. Mini BGD outperforms SGD only if we have a good vectorized implementation. Sum over 10 examples can be performed in a more vectorized way allowing to partially parallelize our computation over the 10 examples, which is not quite possible in case of SGD which looks at just 1 example at a time.

One disadvantage of mini BGD is that there is now this extra parameter  $b$  which may have to be fiddled with that may take time but in a good vectorized implementation mini BGD runs even faster than SGD.

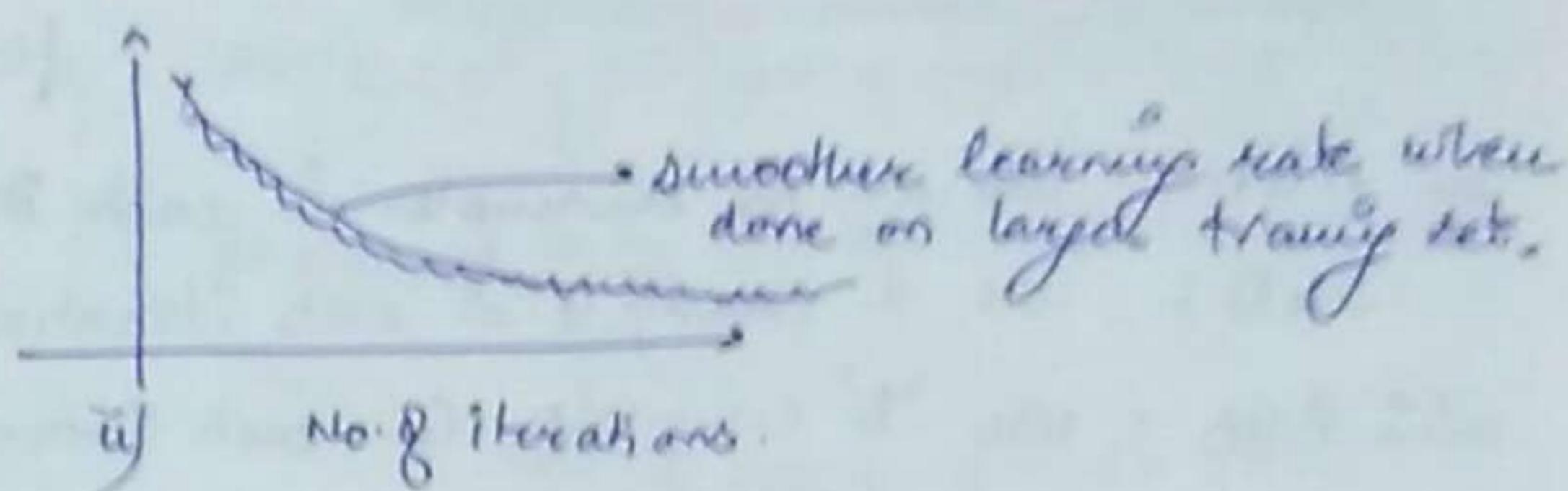
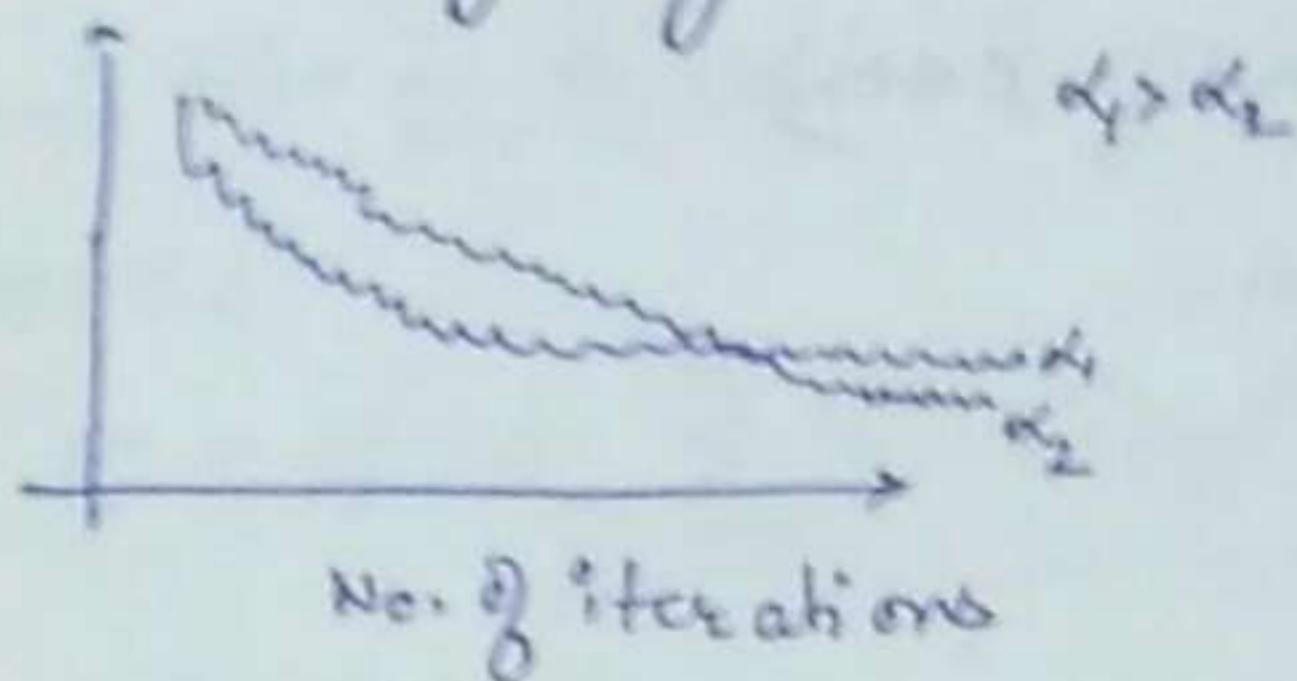
Su BGD,  $b = m$ .

⇒ SGD convergence: to make sure that it's completely debugged and is converging okay, and how to tune  $\alpha$  with SGD.

steps:

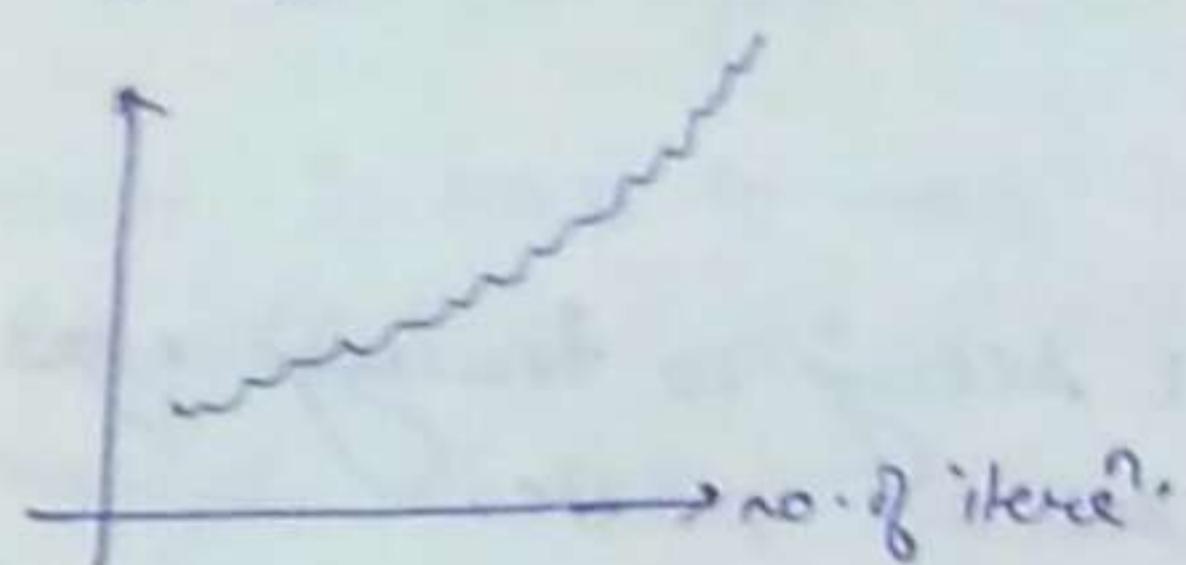
- i) during learning, compute  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$  before updating  $\theta$  using  $(x^{(i)}, y^{(i)})$ .
- ii) every 1000 iterations (say), plot  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$  averaged over last 1000 e.g.s

processed by algorithm.

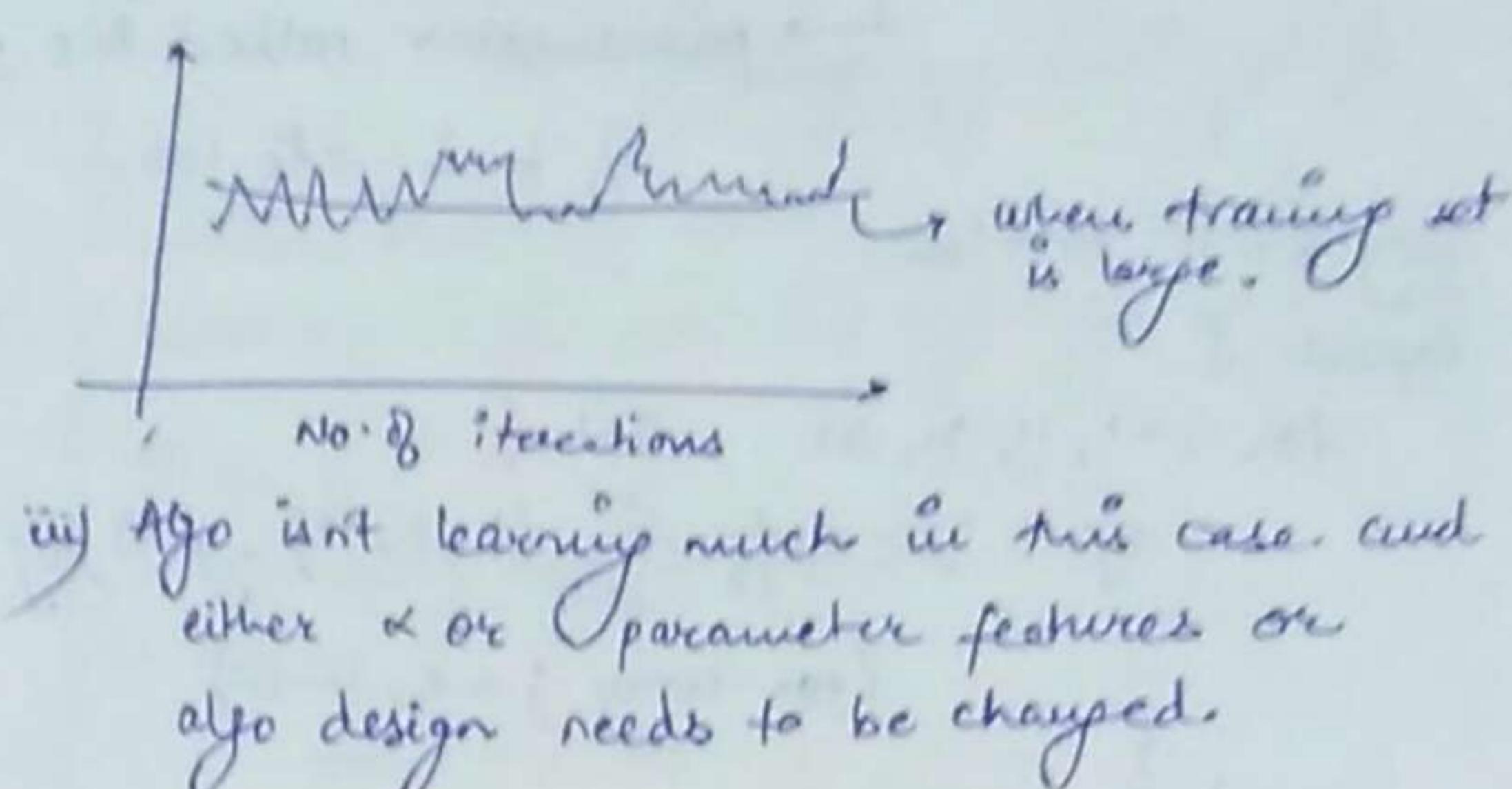


)( Little noisy bcs they are averaged over just a 1000 exgs.) They may not ↓ at every iteration.

Smaller  $\alpha$  results in smaller oscillations and so might give slightly better results.



ii) algo is diverging and smaller  $\alpha$  needs to come.



So, if plot seems noisy ↑ no. of training examples, if learning is not proper change  $\alpha$ .

Note SGD doesn't really converge. It meanders around the global minimum and we get a parameter value that is hopefully close to the global \* but not exact at the global minimum.

Learning rate  $\alpha$  is typically kept constant in SGD and we can slowly ↓  $\alpha$  over time if we want  $\theta$  to converge. but e.g.  $\alpha = \frac{\text{const 1}}{\text{Iteration Number} + \text{const 2}}$

Doing this might lead SGD to exactly converge at the global minimum. People don't prefer doing this much as deciding upon const 1 and const 2 takes a lot of time. The above formula makes sense bcoz as algo runs, the iter? no. becomes large, so  $\alpha$  will slowly become small. One can end up with a slightly better hypothesis if one sets  $\alpha$  slowly to zero.

Though picking very small  $\alpha$  will slow down learning.

Note this method does not require scanning over the entire training set periodically to compute  $J$  on the entire training set, but it looks at just 1000 examples or so.

9/6/18

continuous learning algo which easily adapts.

→ Online learning algo: Many big companies use diff. online learning versions to learn from the flood of users. But the algo there is unlimited data available and there is no need to look at the same data more than once. This algo can adapt to changing user preferences.

e.g. Product search

User searches for "Android phone 1080p camera".

Have 100 phones in store, will return 10 results.

→ features of phone, how many words in user query match name of phone, how many words in query match "descrip" of phone, etc.

$y = 1$  if user clicks on link -  $y = 0$  otherwise. Learn  $p(y=1|x; \theta)$

probability that user will click on a link of a particular phone.

Predicted CLICK THROUGH RATE (CTR)

This CTR can help to show user the 10 phones that they are most likely to click on.

→ choosing special offers to show user, customized "list" of news article, product recommendation. Collaborative filtering can also be used to improve CTR.

There is no need to store the training data and save the learned model rather our algo can continuously evolve with the stream of data coming.

→ Map reduce and data parallelism: Some ML problems are too big to be run on one machine. So, we use

map-reduce approach for large scale ML. In map reduce, we split the training dataset into different subsets and send it to different machines.

$$\text{Machine 1: } \text{temp}_j^{(1)} = \sum_{i=1}^{100} (\text{h}_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

$$\text{Machine 2: } \text{temp}_j^{(2)} = \sum_{i=1}^{100} (\text{h}_{\theta}(x^{(i)}) - y^{(i)}) \cdot w_j^{(i)} \quad \dots \text{and similarly for all other machines upto } m \text{ examples.}$$

Now, each machine has only 1 quarter of the load and could do it about four times as fast.

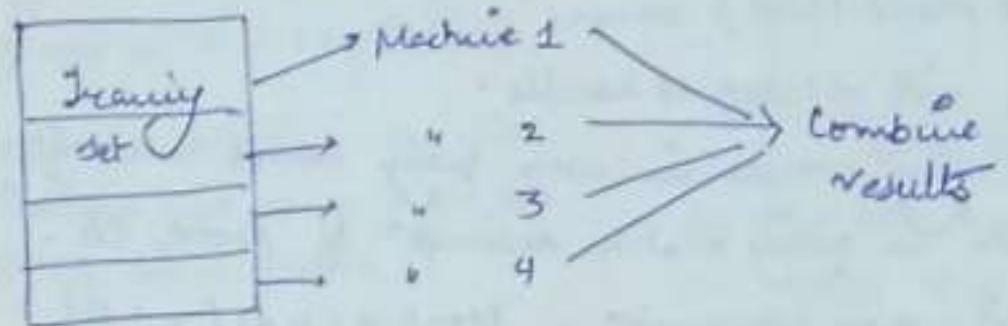
Then all the 'temp' variables are collected and sent to a master server and update parameter  $\theta_j = \theta_j - \alpha \frac{1}{m} (\text{temp}_j^{(1)} + \text{temp}_j^{(2)} + \dots + \text{temp}_j^{(k)})$

(for  $j = 0$  to  $n$ )

$\downarrow$

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

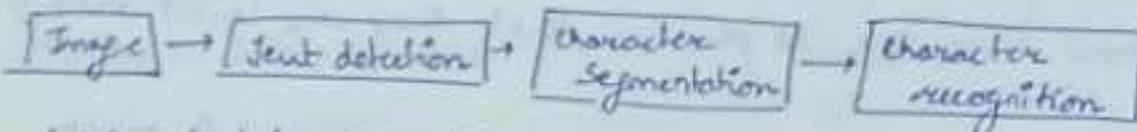
### Map reduce



[giving a 4x speed as 4 computers are in working]

Nok Map reduce is a good candidate when learning algo of ours can be expressed as computing sums of functions over the training set.

### Photo OCR pipeline:



### → Artificial data synthesis:

Creating new data from scratch.

amplify the existing training set  
(using artificial distortions  
on the tnt dataset.)

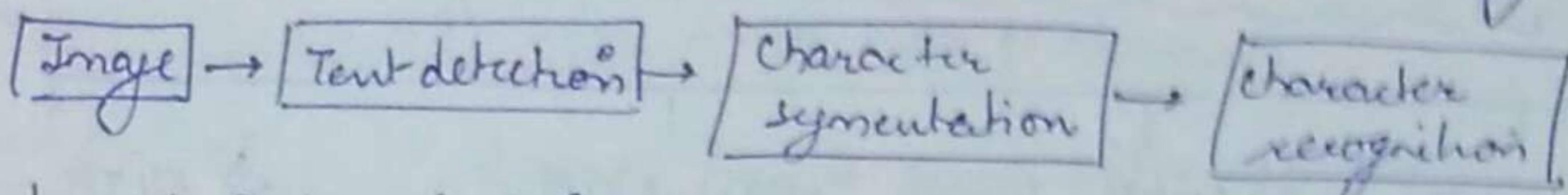
{ Distortion introduced should be representation of the type of noise/distortion in the test set. } Adding meaningless noise is useless.

Nok make sure to have a low bias classifier before expanding the training data. Keep trip size no. of features in row until we have a low bias classifier.

"Crowd source" is the method of labelling data by human labour.

"Celery analysis" : to give guidance on what parts of the pipeline might be best use of your time to work on.

e.g. Estimating the errors due to each component (ceilng analysis).



What part of the pipeline should you spend the most time trying to improve?

Component	Accuracy
Overall system	72%
Text detection	89% ↓ 17% ↑
Char segmenta <sup>n</sup>	90% ↓ only 1% ↑
Char recogni <sup>n</sup>	100% ↓ 10% ↑

what we are doing is first we checked the accuracy for overall system. then we checked the accuracy of the system by having 100% correct text detection (by manually segmenting the text) and found that when this segment of pipeline achieved 100% accuracy, the final results improved by 17%. Hence, it would be worth spending time on improving the accuracy of this part of our segment rather than the character segmentation part bcoz even if it does with 100% perfect, the overall accuracy gain is just 1%.

'Ceiling analysis' also helps us find the max theoretical accuracy that can be achieved if all the segments of the pipeline worked perfectly.