

Machine Learning - Andrew Ng.

ML

$$\mathbf{x} = \{x_0^1, x_1^1, x_2^1\} \\ x_0^1 = 1, x_1^1 = 2, x_2^1 = 3$$

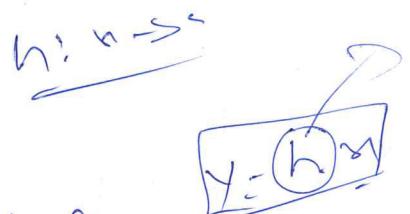
① $\mathbb{R}^n \rightarrow \mathbb{R}$

x : set of inputs of size n

$$x = (x^n)$$

y : set of outputs

$$y = y^n$$



$\rightarrow h: \mathbb{R}^n \rightarrow \mathbb{R}$, when ~~\mathbb{R}^n~~ , $y \in \mathbb{R}$

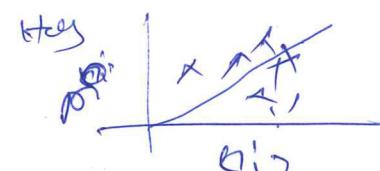
\rightarrow when $n = \mathbb{R}^m$ and $y = \mathbb{R}^p$, the simplest choice is to approximate y as a linear function of x

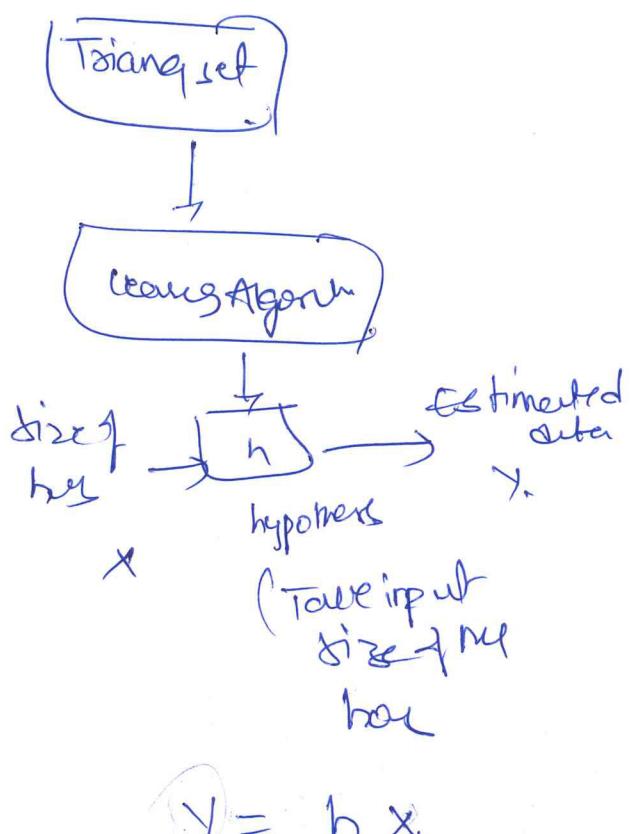
$$h_{\text{linear}} = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n.$$

x

size \mathbb{R}^n choice \mathbb{R}^p (x)

$x(1)$	$y(1)$
$x(2)$	$y(2)$
$x(3)$	$y(3)$





h = a function that maps x 's to y 's,

$$h(x) = \theta_0 + \theta_1 x.$$

(one variable)

\Rightarrow one variable

Linear

h

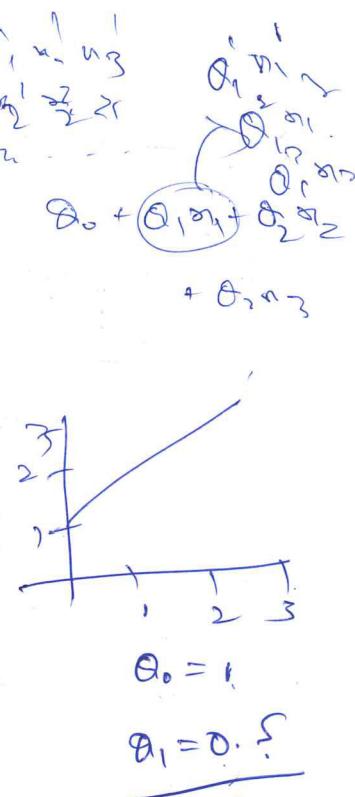
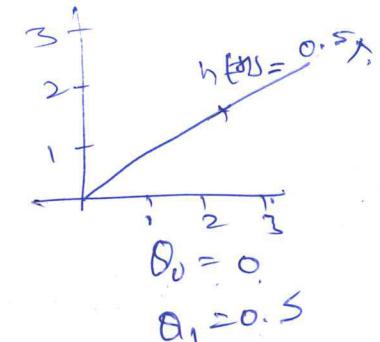
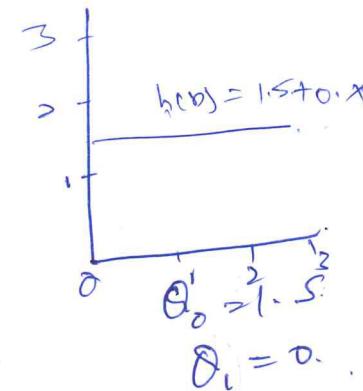
cost function): How to fit best possible straight line to our data.

$$\begin{matrix} x_1 & y_1 \\ x_2 & y_2 \\ x_3 & y_3 \\ x_4 & y_4 \\ x_5 & y_5 \\ x_6 & y_6 \\ x_7 & y_7 \\ x_8 & y_8 \\ x_9 & y_9 \end{matrix}$$

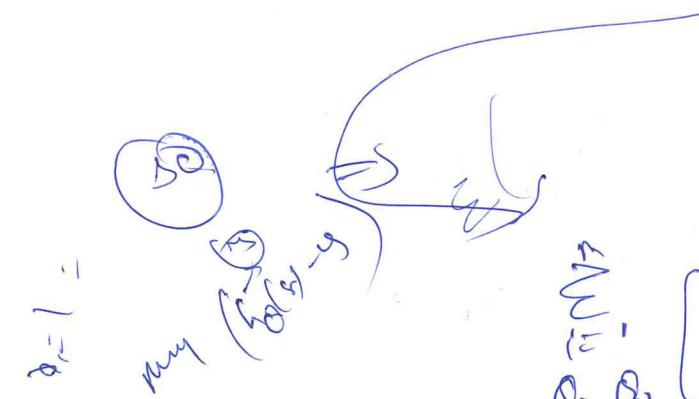
$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

θ 's = parameters

How to choose θ_0 and θ_1 ,



\Rightarrow so minimize (θ_0, θ_1)



$$\text{so } \sum_{i=1}^m (h_i(\theta) - y_i)^2$$

Difference has to be minimum

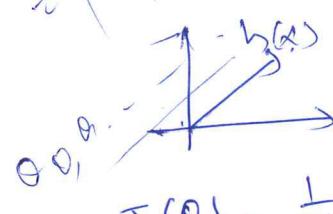
$$\sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2$$

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

minimize θ_0, θ_1 $J(\theta_0, \theta_1) \rightarrow$ square error cost function

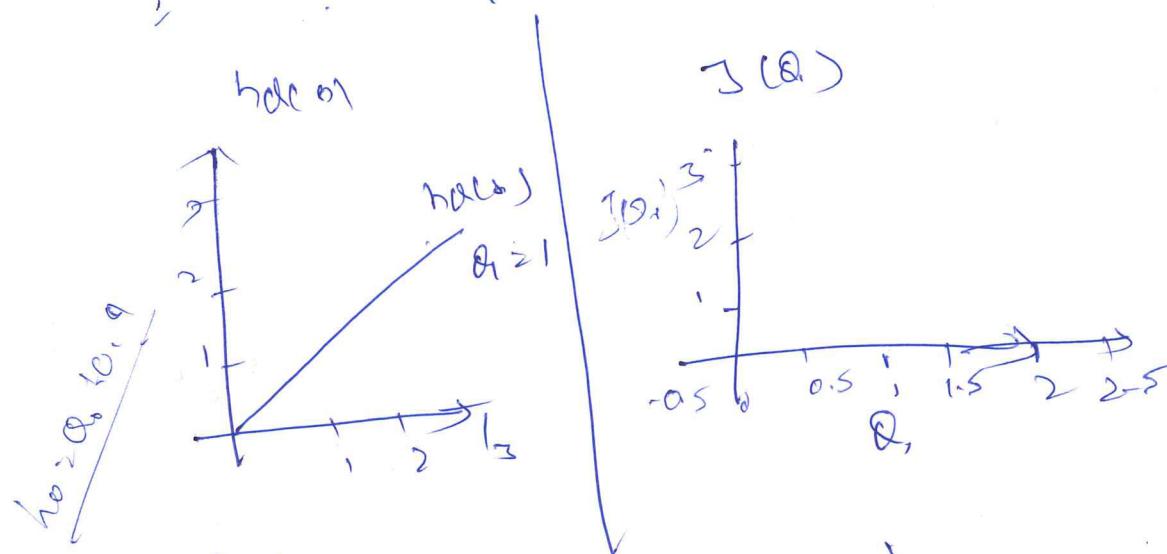
simplified

$$h_\theta(\theta_0) = \theta_1 x \therefore \theta_0 = 0.$$



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

hold on



$$\theta_1 = 1$$

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

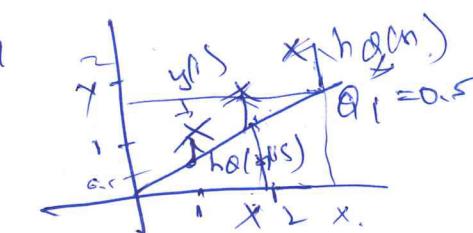
$$= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 =$$

constant for θ_1 and θ_2

$$= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0^2$$

$$J(1) = 0$$

If $\theta_1 = 0.5$



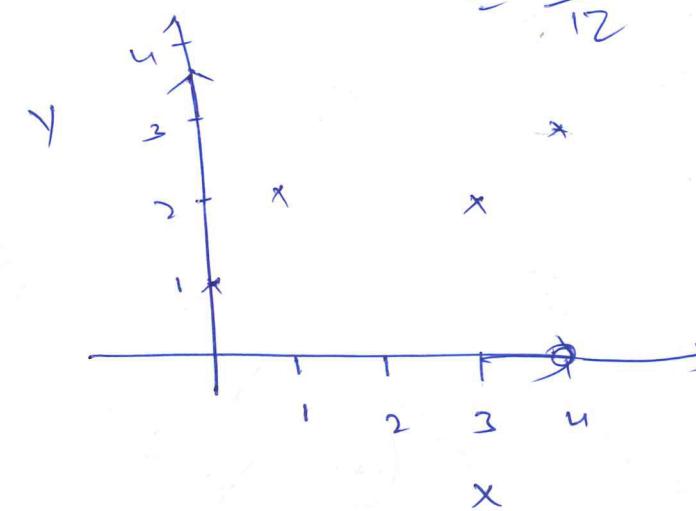
$$J(0.5) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$= m = 3$$

$$= \frac{1}{3} [(0.5)^2 + (1)^2 + (1.5)^2]$$

$$= \frac{1}{3} \left[\frac{7}{2} \right]$$

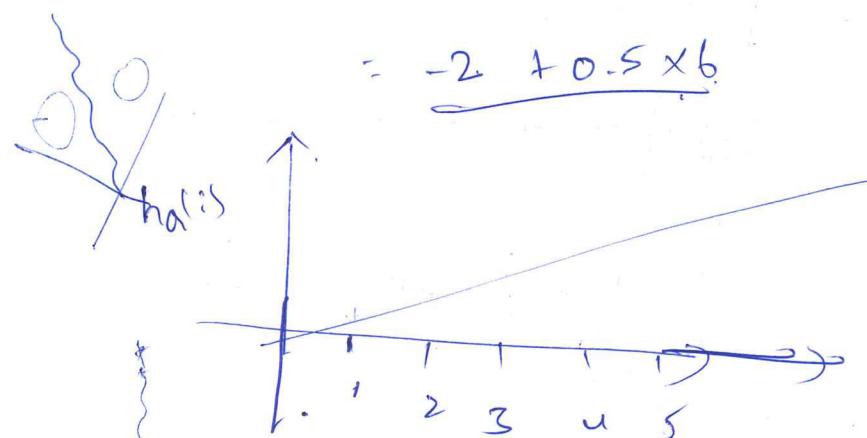
$$= \frac{7}{12}$$



-1780.0 + 11 (S₂₀₋₉)

$$\cancel{h(0)} = \frac{1}{\cancel{60}} [$$

$$\text{Holes} = Q_1 + Q_2 \times 6. \quad h_{D0} = Q_1 + Q_2 \times b$$



$$h(0) \approx -1.5 \text{ N/m}$$

$$h'(2) = -2 + 0 \cdot S(2) = -2$$

$$h'(2) = -2 + 0 \cdot S(2) = -2$$

$$h'(2) = -2 + 0 \cdot S(2) = -2$$

$$\begin{aligned} h(2) &= -2 + 0 \\ h(3) &= -2 + 0.5(3) = 0.5 \\ h(0.5) &= -2 + 0.5(0.5) = 0 \end{aligned}$$

$$h_2(t) = -2 + 0.5t \Rightarrow$$

$$\cancel{b_0(2)} \rightarrow 2 + 0^{-5}$$

$$T = 2 + 0 \cdot s(3)$$

$$\text{b) } h_0(1) = -2+0 \cdot 1 \Rightarrow h_0(1) \Rightarrow h_0(0) = -2$$

$$h_{\alpha}(x^2) = -2 \cos(\alpha) \Rightarrow h_{\alpha}''(0) = -2 \cos(\alpha)$$

$$h_0(3) = -2 \text{to.} \leq (u) = h_0(u) = -\infty$$

1. $\log(2^n)$

15

$$\begin{aligned} & \text{no}(s) \xrightarrow{\sim} \text{no}(s') \\ & \text{no}(s') \xrightarrow{\sim} \text{no}(s'') \end{aligned}$$

$\Theta \Rightarrow$

1

$$3 \cdot z^1$$

$$(0 + 0.5 - 0.5)^2 =$$

$$(0+1-1)^2$$

(0 + 2)

$\text{f}(\text{f}^{-1}(x)) = x$

10

$$2 + 0 \cdot (3) =$$

$$h_0(2) = 0 + 1 \leftarrow (2)$$

$$= 0 + \textcircled{3} \quad \textcircled{1}$$

matrices :- Wk 4 - 2

$$\begin{bmatrix} 1 & 3 \\ 2 & 4 \\ 0 & 5 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 3 \\ 2 & 4 \\ 0 & 5 \end{bmatrix}, B = \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix}$$

$$AB = \begin{bmatrix} 1 & 3 \\ 2 & 4 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} 7 & 9 \\ 14 & 12 \\ 0 & 15 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & -4 \\ 2 & 1 \end{bmatrix}, B = \begin{bmatrix} 0 & 3 \\ 5 & 8 \end{bmatrix}$$

$$A + B = \begin{bmatrix} 1 & -2 \\ -4 & 1 \end{bmatrix}$$

$$B = \begin{bmatrix} 0 & 3 \\ 5 & 8 \end{bmatrix} = A + B = \begin{bmatrix} 1 & 3 \\ -1 & 9 \end{bmatrix}$$

$$v = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix}, w = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix}$$

$$vw^T = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \begin{bmatrix} 2 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 2+6+4 \\ 6+6+2 \\ 2+2+1 \end{bmatrix} = \begin{bmatrix} 12 \\ 14 \\ 5 \end{bmatrix}$$

$$(A)(B)$$

$$v = \begin{bmatrix} -3 & 3 \\ 4 & 3 \end{bmatrix}, w = \begin{bmatrix} 3 & 1 \\ 2 & 5 \end{bmatrix}$$

$$v = \begin{bmatrix} 3 & 1 \\ 2 & 5 \end{bmatrix} = \begin{bmatrix} -9+15+12 \\ -6+5+10 \end{bmatrix}$$

* Linear regression - multivariate.

~~Im~~ \rightarrow $\begin{bmatrix} n & 1 & y \end{bmatrix} \rightarrow$ one
one \Rightarrow single.

In	Fr	In	Out
x_1	x_{12}	x_{13}	y
x_1	x_{21}	x_{22}	y
x_1	x_{22}	x_{23}	y
x_1	x_{31}	x_{32}	y
x_1	x_{32}	x_{33}	y

n = number of feature.

$\xrightarrow{\text{in}}$ $\xrightarrow{\text{out}}$

$x_{ij} = x_i \xrightarrow{\text{in}} \xrightarrow{\text{out}}$

$x_{ij} = 1$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$h_\theta(x^1) = \theta_0 + \theta_1 x_1^1 + \theta_2 x_2^1 + \theta_3 x_3^1 + \dots + \theta_n x_n^1$$

$$h_\theta(x^1) = \theta_0 + \theta_1 x_1^1 + \theta_2 x_2^1 + \dots$$

$$h_\theta(x^2) = \theta_0 + \theta_1 (x_1^2) + \theta_2 (x_2^2) + \dots$$

$$\begin{bmatrix} h_{\theta}(x^1) \\ h_{\theta}(x^2) \\ h_{\theta}(x^3) \end{bmatrix} = \begin{bmatrix} x_0 & x_1 & x_2 & x_3 \\ x_0^2 & x_1^2 & x_2^2 & x_3^2 \\ x_0^3 & x_1^3 & x_2^3 & x_3^3 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$= \theta^T \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$= [\theta_0, \theta_1, \theta_2] \begin{bmatrix} x_0 & x_1 & x_2 & x_3 \\ x_0^2 & x_1^2 & x_2^2 & x_3^2 \\ x_0^3 & x_1^3 & x_2^3 & x_3^3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}, x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

~~Octave~~

$$a = [1 \ 2 \ 3 \ 5 \ 6]$$

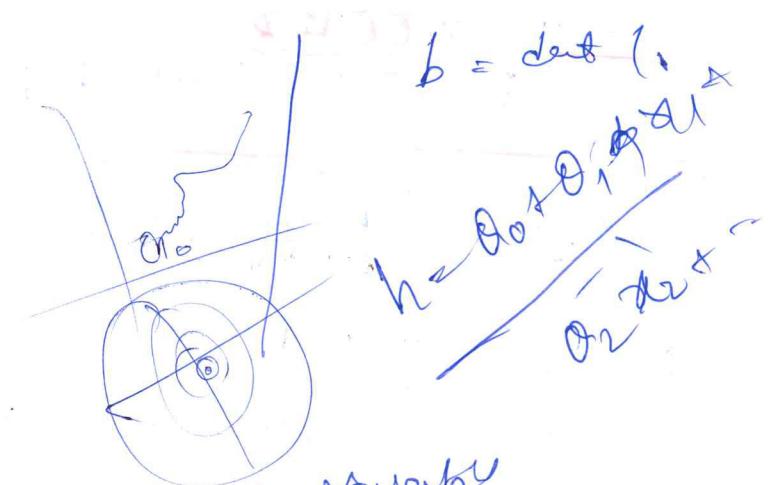
$$a = \begin{matrix} 1 & 2 & 3 \\ 5 & 6 \end{matrix}$$

$$b = a(1,1)$$

$$b = \frac{1}{4}$$

$$b = \alpha(i, 3)$$

$$\begin{matrix} 3 \\ 6 \\ b = (1, 1) \\ \underline{1 \ 2 \ 3} \end{matrix}$$



Gradient descent for multiple variables

$$\begin{cases} \theta_0: \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)} \\ \theta_1: \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)} \\ \theta_2: \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)} \end{cases}$$

~~Update θ_0, θ_1~~

~~Update θ_1, θ_2~~

~~Update θ_2~~

~~Update θ_0~~

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

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

~~$= \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$~~

Update augman

Gradient descent method

One variable Gradient descent method

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$\frac{\partial J(\theta)}{\partial \theta}$

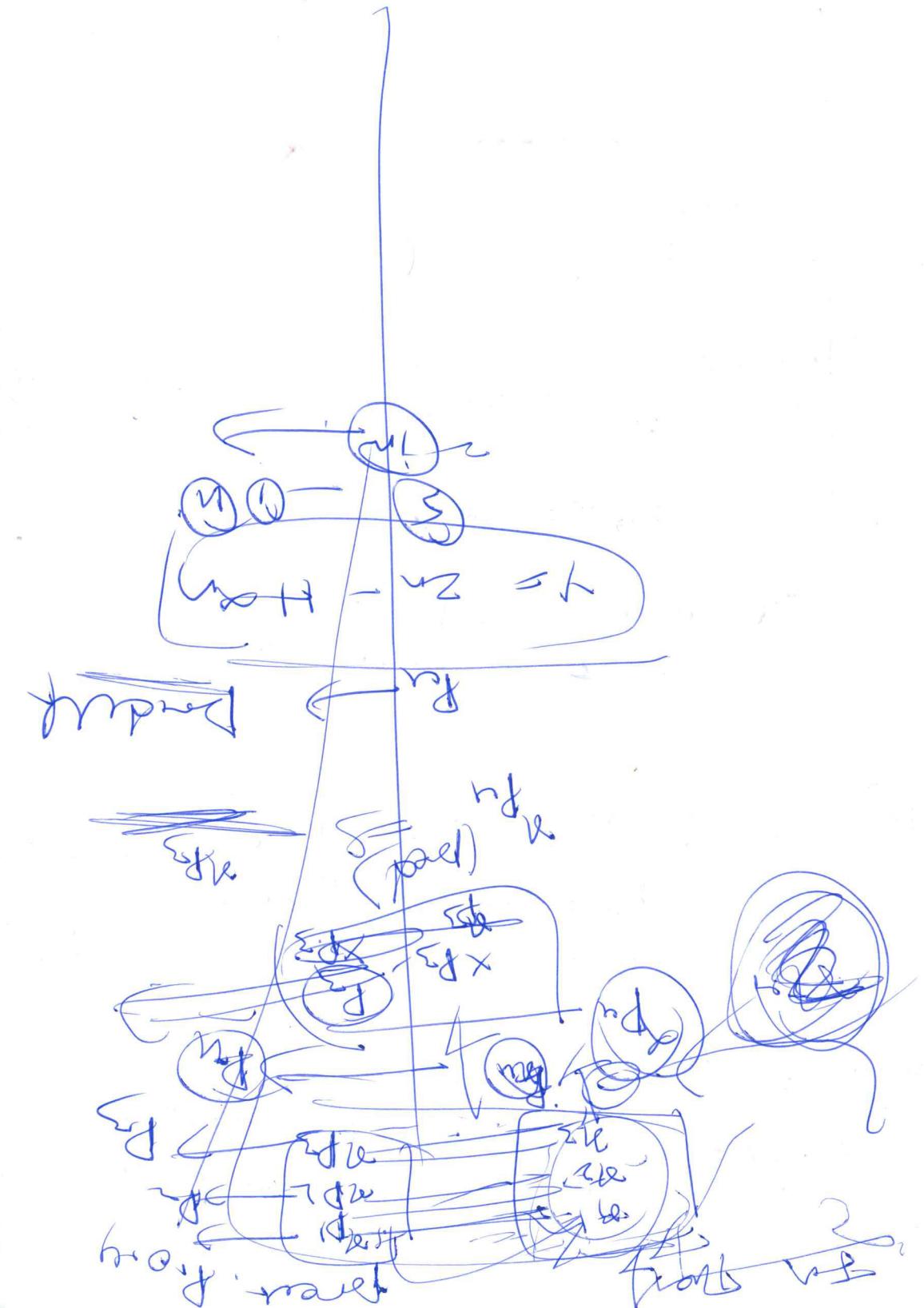
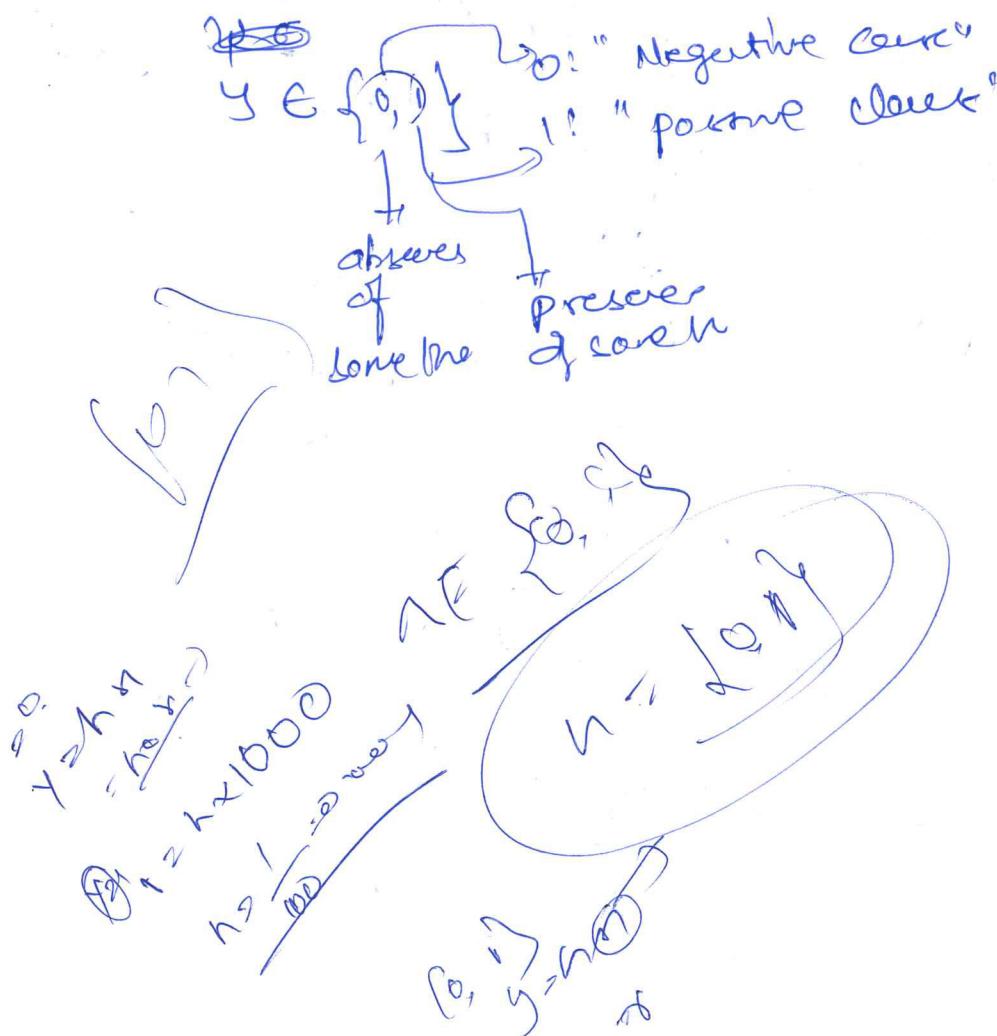
Week 3: WEEK 3: CLASSIFICATION AND Representation

Logistic Regression

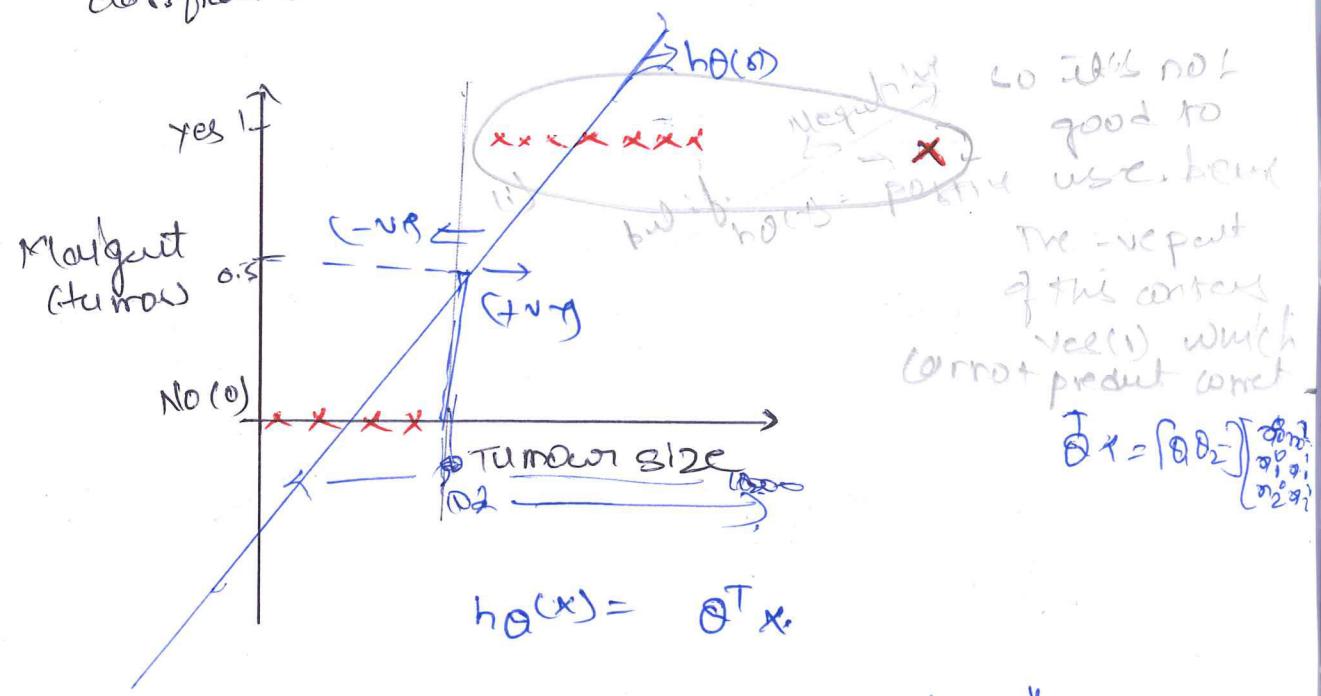
Classification :-
Email is - spam / notspam?

Error! :- spam | notspam?

online : foodbev (yes / no))



mostly classification: $y \in \{0, 1, 2, 3\}$ =
classification



If $\text{hours} \geq 0.5$, predict " $y=1$ "
 $\text{hours} < 0.5$, predict " $y=0$ "

- For Classification: $\gamma = 0, 0.5, 1$
- $h_\theta(x)$ can be > 1 or < 0
- { So error
 Not good
- So, we go to

Logistic Regression: $\Phi \circ \in \text{hol}(\mathcal{S})$
is also rate of classifier

Logistic Regression: Hypothesis Regression.

Want $\Sigma \leq \text{halos} \leq 1$

$$\text{so } \underline{\underline{h_{\theta}(x)}} = \Theta x.$$

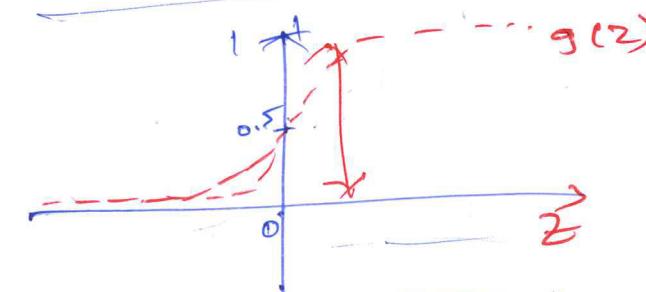
$$\hookrightarrow h_{\theta}(z) = g(\theta T_z) \\ = g(z)$$

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

$$= \frac{1}{1 + e^{-z}}$$

$$h_{\theta}^{(0)} = \frac{1}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{-2}}$$

Sigmoid function or Logistic function



Interpretation of Hypothesis Output

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x

$$\text{Ex: } \text{If } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorsize} \end{bmatrix}$$

$$y=1$$

$h_{\theta}(x) = 0.7$, probability of $y=1$ is 0.7

(*) 70% chance of tumor being malignant

$h_{\theta}(x) = P(y=1/x; \theta) \rightarrow$ "probability that $y=1$, given x , parameterized by θ "

→ As this is a classification

$y \geq 0$ or 1.

so we can calculate $y \geq 0$ for given

~~$P(y=0/x; \theta) = 0.3,$~~

as $P(y=1/x; \theta) + P(y=0/x; \theta) = 1$

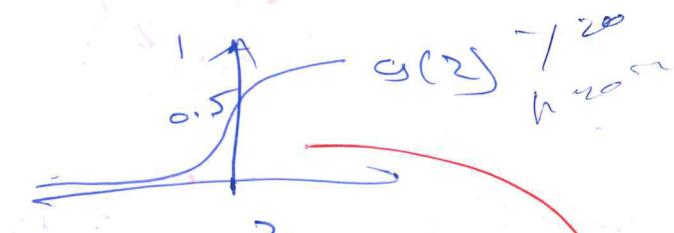
~~$P(y=0/x; \theta) = 1 - P(y=1/x; \theta) = 1 - 0.7 = 0.3$~~

Decision Boundary

$$\log(h_{\theta}(x))$$

$$h_{\theta}(x) = g(\theta^T x) = P(y=1/x; \theta)$$

when it makes $y=1$ and $y=0$ and to better understand the hypothesis (particularly when there are many features).



Suppose predict " $y=1$ " if $h_{\theta}(x) \geq 0.5$

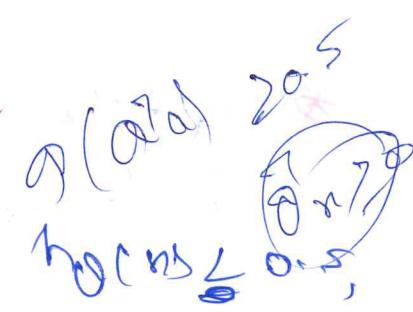
predict " $y=0$ " if $h_{\theta}(x) < 0.5$

→ To understand better, we need when

~~$y \geq 0$ given $h_{\theta}(x) \geq 0.5$~~

when $g(z) \geq 0.5$
when $z \geq 0$

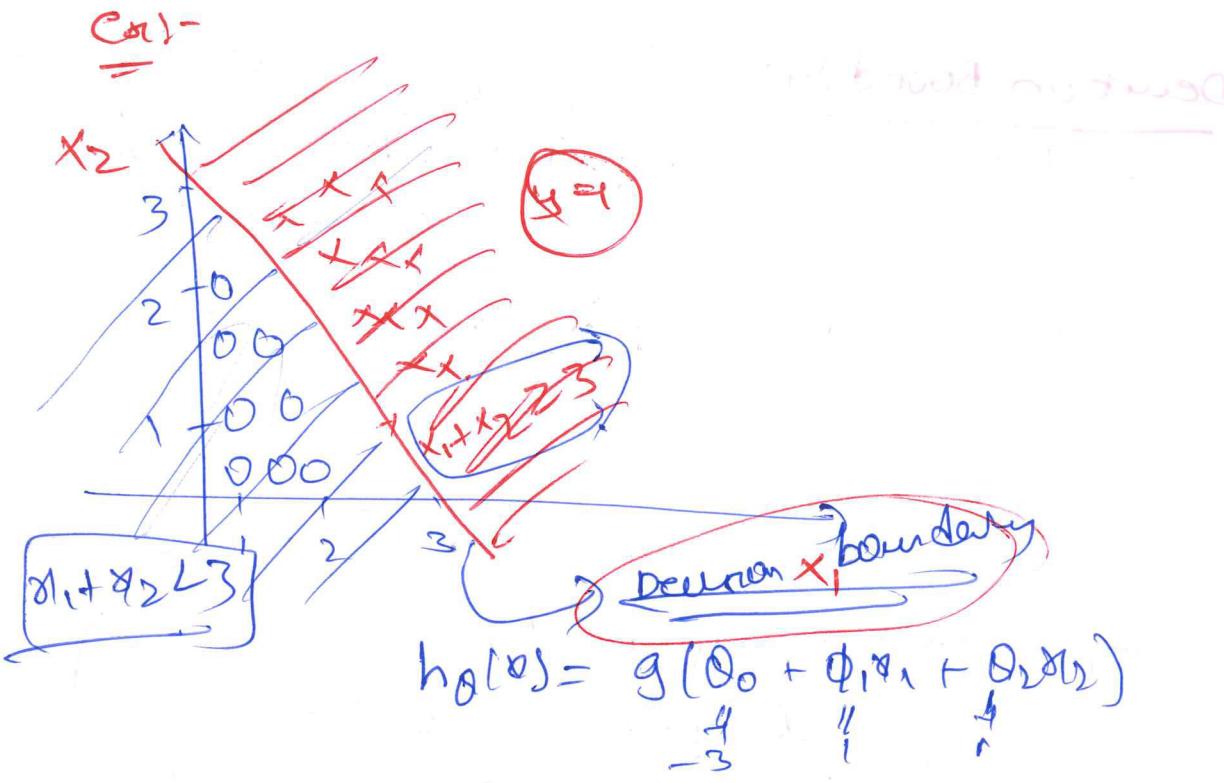
so $h_{\theta}(x) = g(\theta^T x) \geq 0.5$, when $z \geq 0$.



$$\theta^T x \geq 0$$

→ Similarly $h_{\theta}(x) \leq 0.5$, when

$$\theta^T x < 0$$



$y=1, \text{ if } (-3 + \theta_1 x_1 + \theta_2 x_2 \geq 0)$

$y=1, \text{ if } x_1 + x_2 \geq 3$

~~$x_1 + x_2$~~

$y=0 \rightarrow \text{if } (-3 + \theta_1 x_1 + \theta_2 x_2 \leq 0)$

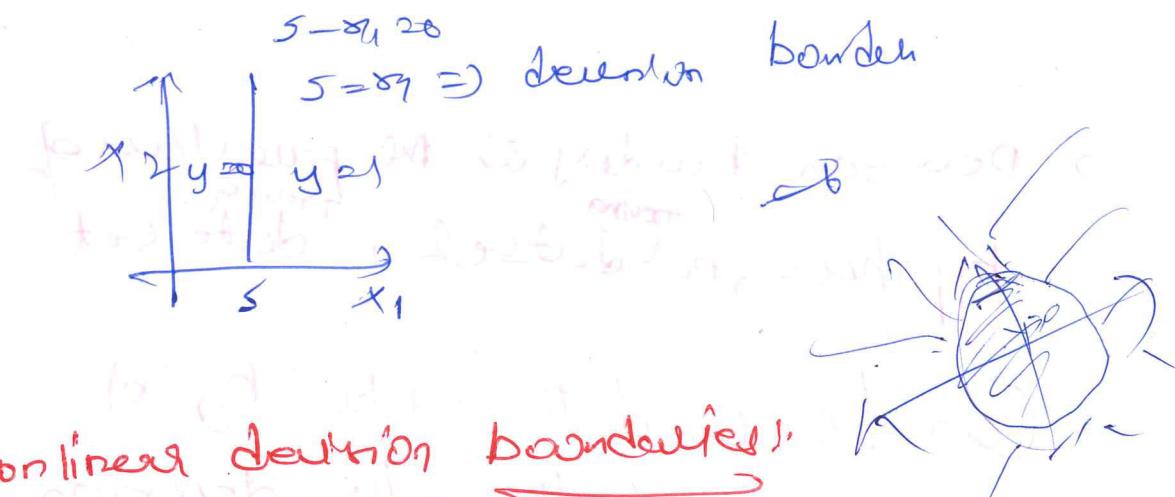
$x_1 + x_2 < 3 \Rightarrow y=0.$

- These are used,
- Decision boundary is a property of $h_0(x)$ involving $\theta_0, \theta_1, \theta_2$.
- These are in general property of $h_0(x)$ not dataset

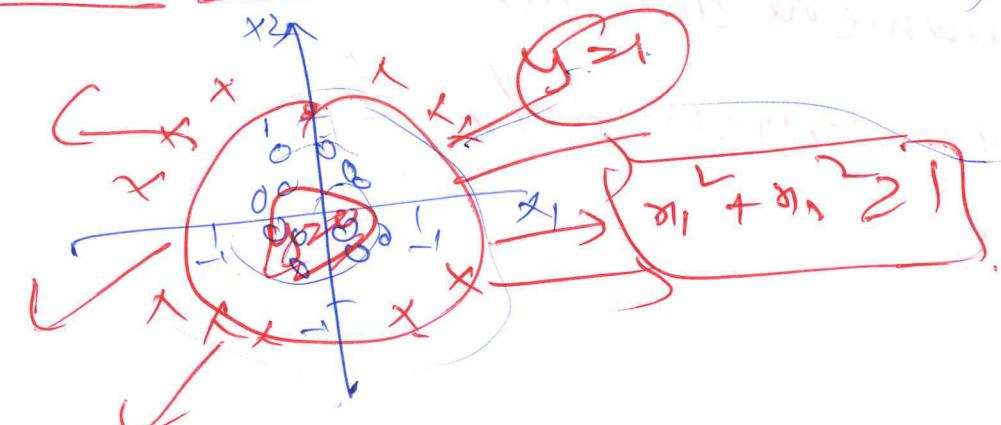
cont-

$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta_0 = -3, \theta_1 = 1, \theta_2 = 1$$



Non linear decision boundaries



$$h_{\theta}(\mathbf{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \frac{\theta_3}{2} x_1^2 + \frac{\theta_4}{4} x_2^2)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$$

Predict y' if $-x_1 + x_1^2 + x_2^2 \geq 0$

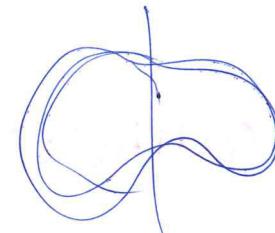
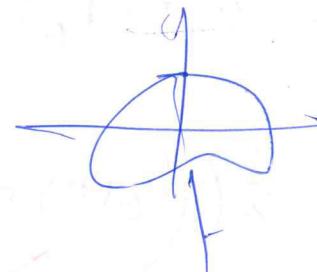
$$x_1^2 + x_2^2 \geq 1 \Rightarrow \begin{cases} x_1 \geq 0, x_2 \geq 1 \\ x_2 \geq 0, x_1 \geq 1 \end{cases}$$

\rightarrow Decision boundary is the ~~perpendiculars of~~
~~hypotenuse, not vertical~~ ~~training~~ ~~test~~ ~~dataset~~

Can be used to fit the $h_{\theta}(\mathbf{x})$
parameters and then the decision
boundary is defined.

* more complex decision boundary

$$h_{\theta}(\mathbf{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \dots)$$



Logistic Regression: Cost Function: How to

auto fit parameter θ to ~~to fit~~ ~~to training~~
but

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

measure no. of $\{x^{(i)}\}_{i=1}^{n_0}$, $n_0 = 1, 4, 5, 9, 14$

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

How to choose θ ?

cost function

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

Logistic Regressn.

$$\text{cost}(h_{\theta}(x), y)$$

$$\rightarrow \text{lost} = 0, \text{ if } y = 1, h_{\theta}(y) = 1$$

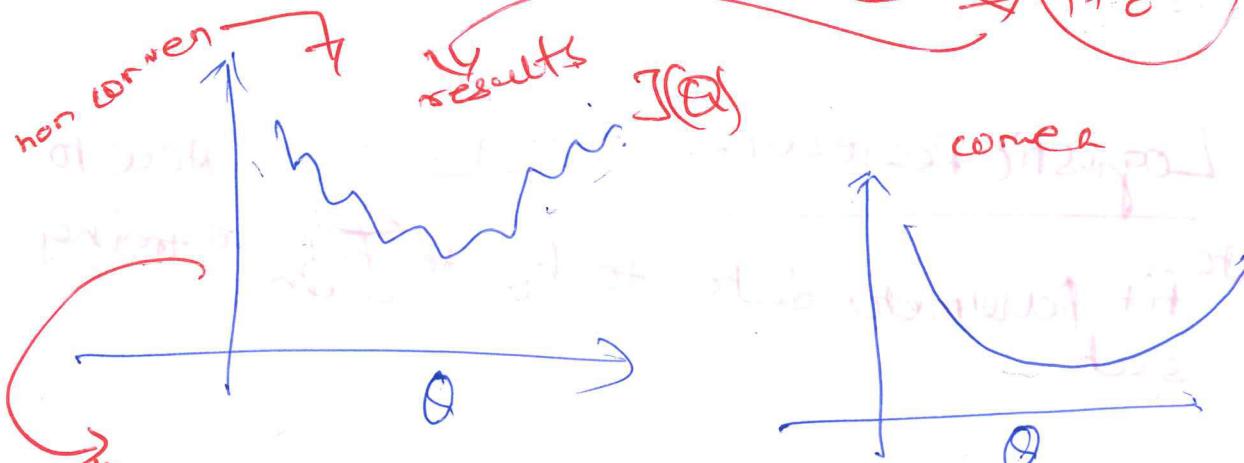
But as $h_{\theta}(y) \rightarrow 0$

cost $\rightarrow \infty$.

If $h_{\theta}(y) = 0$,

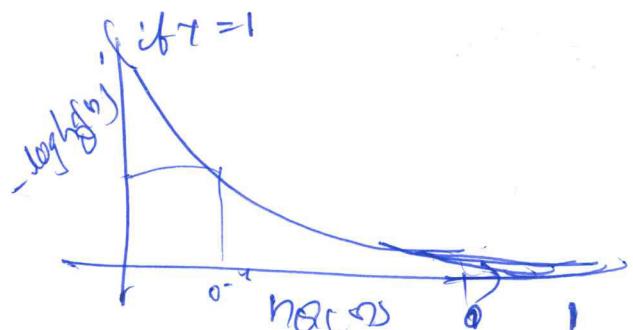
(Predictor, $P(y=1|x, \theta) = 0$) but $y=1$
we'll perceive linear algorithm very by
very large cost.

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

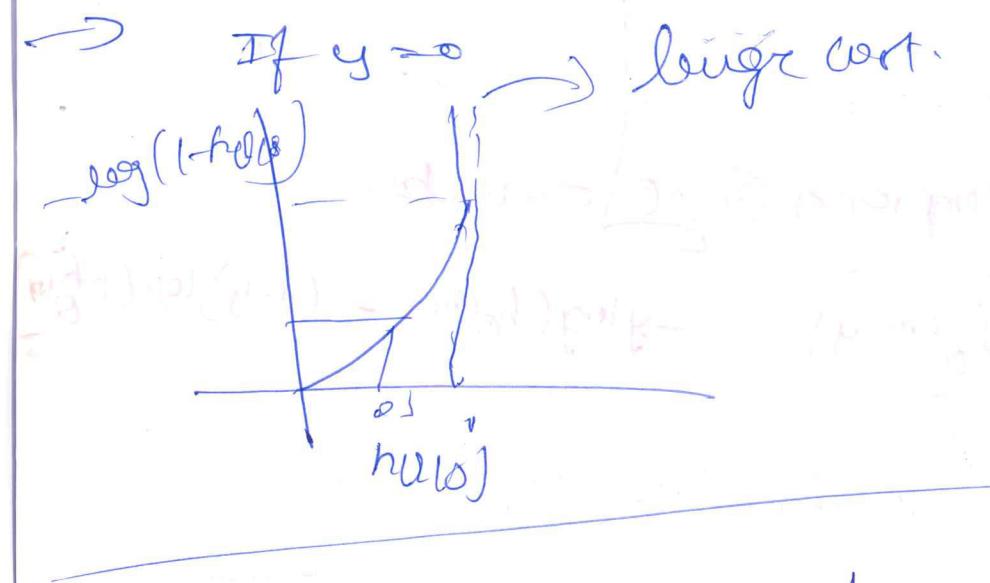


* SO we cannot use the same cost function for logistic, instead to avoid many.
Logistic regression cost function so use

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



$$\log \frac{1}{1+h_{\theta}(x)}$$



Sigm

$$-\log(h_{\theta}(x)) \quad y=1$$

$$-\log(1-h_{\theta}(x)) \quad y=0$$

$$-\log\left(\frac{1}{1+h_{\theta}(x)}\right)$$

$$-\log\left(\frac{1}{1+h_{\theta}(x)}\right) \quad y=0$$

$$-\log\left(\frac{h_{\theta}(x)}{1-h_{\theta}(x)}\right) \quad y=1$$

$$-\mu(\text{correct}-1)$$

Simplified cost function and Gradient descent

descent:

Logistic regression (P2)

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

$$\rightarrow \text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)), & y=1 \\ -\log(1-h_{\theta}(x)) & y=0 \end{cases}$$

Note $y=0, 1$ cases

so the cost function can be ordered to

comprises ② in ① criteria

$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$\text{If } y=1; \text{ Cost}(h_{\theta}(x), y) = -\log(h_{\theta}(x))$$

$$\text{If } y=0; \text{ Cost}(h_{\theta}(x), y) = -\log(1-h_{\theta}(x))$$

$$\text{So: } J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))]$$

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

To fit parameters θ :

$$\min_{\theta} J(\theta) \rightarrow \text{cost} \theta$$

① To make a prediction given new x :

$$\text{Output } h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \quad P(y=1/x; \theta)$$

② By Gradient descent

$$\text{Want } \min_{\theta} J(\theta)$$

Repeat {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

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

↑ Iteratively update all θ_j

Logistic Regression: Advanced Optimization

cost function $J(\theta)$ • Want min $J(\theta)$

Given $\theta_j \rightarrow$ code to compute

$$\begin{aligned} -J(\theta) \\ -\frac{\partial}{\partial \theta_j} J(\theta) \end{aligned}$$

Gradient descent

repeat {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

Optimization Algorithms

→ Gradient descent

- conjugate gradient

- BFGS

- LBFGS

} out of slope

Example:

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\theta) = (\theta_1 - s)^2 + (\theta_2 - t)^2$$

$$\min_{\theta} J(\theta)$$

$$\frac{\partial \theta}{\partial \theta_1} J(\theta) = 2(\theta_1 - s)$$

$$\frac{\partial \theta}{\partial \theta_2} J(\theta) = 2(\theta_2 - t)$$

octave example

%%%

$$\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta)$$

for loop

function [Jval, grad] = costFunction(theta)

jVal = [- - code to compute $J(\theta)$]

grad = [- - code to compute derivative
of $J(\theta)$];

end.

Logistic Regression: MULTI CLASSIFICATION (One-vs-all)

Email :- work, friends, Family, Hobby
 $y=1$ $y=2$ $y=3$ $y=4$

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

Medical diagrams! ~~Not~~ Not ill, cold, FLU
 $T_y=1$ $T_y=2$ $T_y=3$

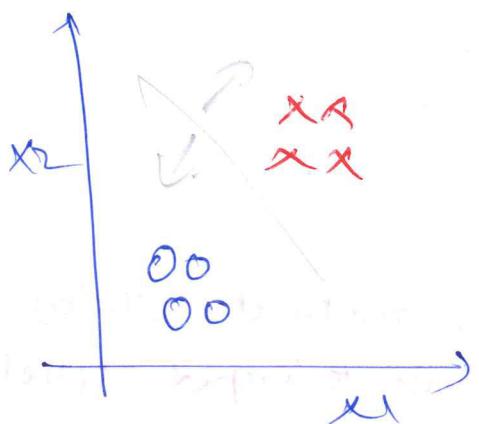
one-on-one clearing has for every class of file
to predict the probability that $y = 1$

Weather : Sunny, cloudy, Rain, Snow

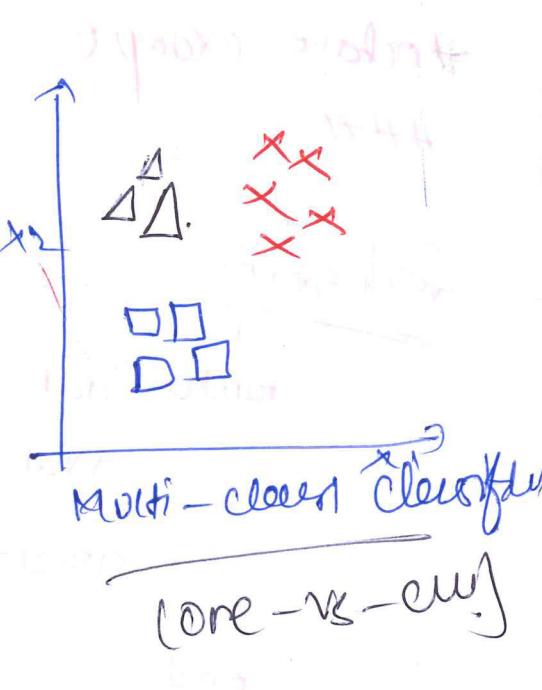
$y=1$ $y=2$ $y=3$ $y=4$

t_0 t_1 t_2 t_3

On a new input x , to make a prediction, pull the
class(es) that maximises $\text{max}_i h_i^{(c)}(x)$.



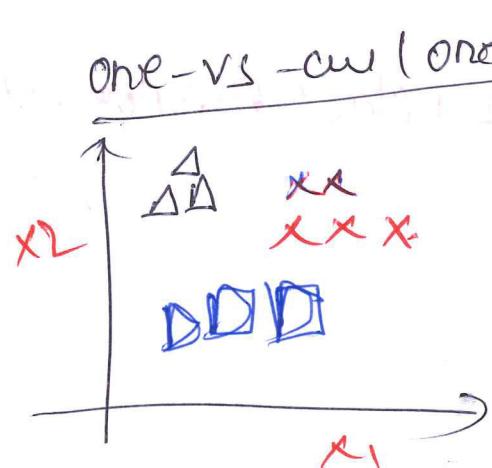
Binary Classification



Class: 1

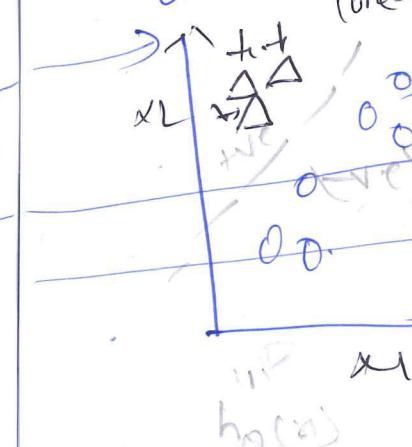
class ②

Class 8



Class 1:
Class 2:
Class 3:

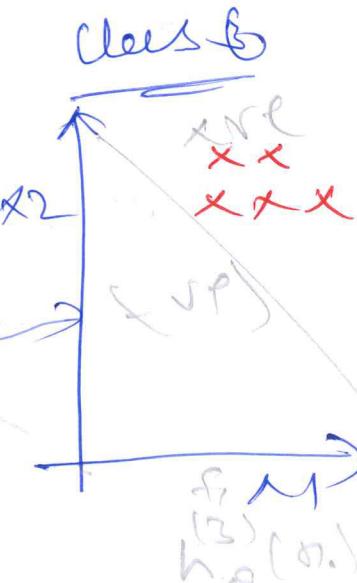
Class: ①



$\begin{array}{r} 00 \\ 0 \\ \hline 2 \end{array}$ $\begin{array}{r} 00 \\ 0 \\ \hline - 5 e \end{array}$

DDD (Ans)

X



ex 2 Q12

x_0	x_1	x_2	y
1	1	0.5	0
1	1	1.5	0
2	1	1	1
3	1	1	0

$$h = \theta^T x$$

$$\begin{bmatrix} h_0 \\ h_1 \\ h_2 \end{bmatrix} = \theta \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 \end{bmatrix} \begin{bmatrix} x_0 & x_1 & x_2 \\ x_0^2 & x_1^2 & x_2^2 \\ x_0 x_1 & x_1 x_2 & x_0 x_2 \end{bmatrix}$$

$$h = \begin{bmatrix} \theta_0 + \theta_1 + 0.5\theta_2 \\ \theta_0 + \theta_1 + 1.5\theta_2 \\ \theta_0 + 2\theta_1 + \theta_2 \\ \theta_0 + 3\theta_1 + 0.5\theta_2 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\cos(h - y) =$$

$$b + x_1 \geq 0$$

$$x_1 \geq 0$$

$$x_1 \leq 7$$

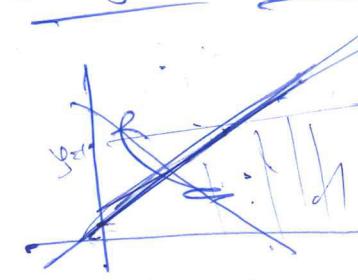
$$x_1 \geq 6$$

$$m \leq 6$$

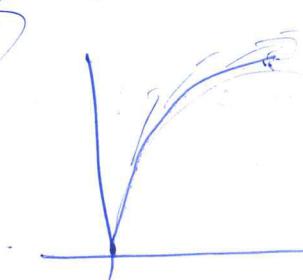
$$x_1 \geq 6$$

$$x_1 \geq 6$$

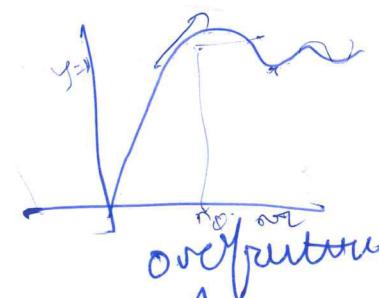
* Different modes of Overfitting



underfit



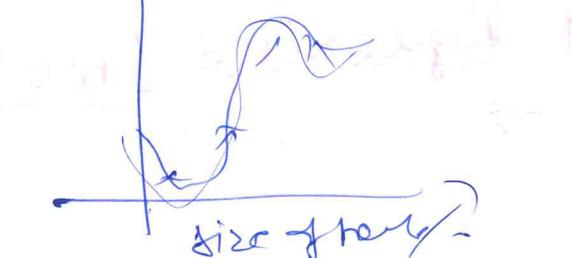
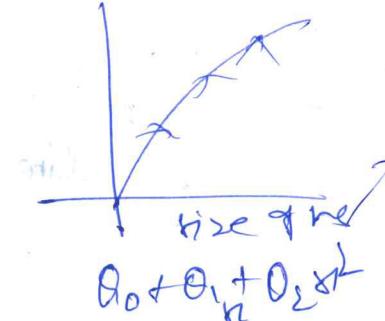
overfit



overfit

Can be avoided by choosing
the no of parameters.

cost function



$$\frac{\partial}{\partial \theta} J(\theta) = 0$$

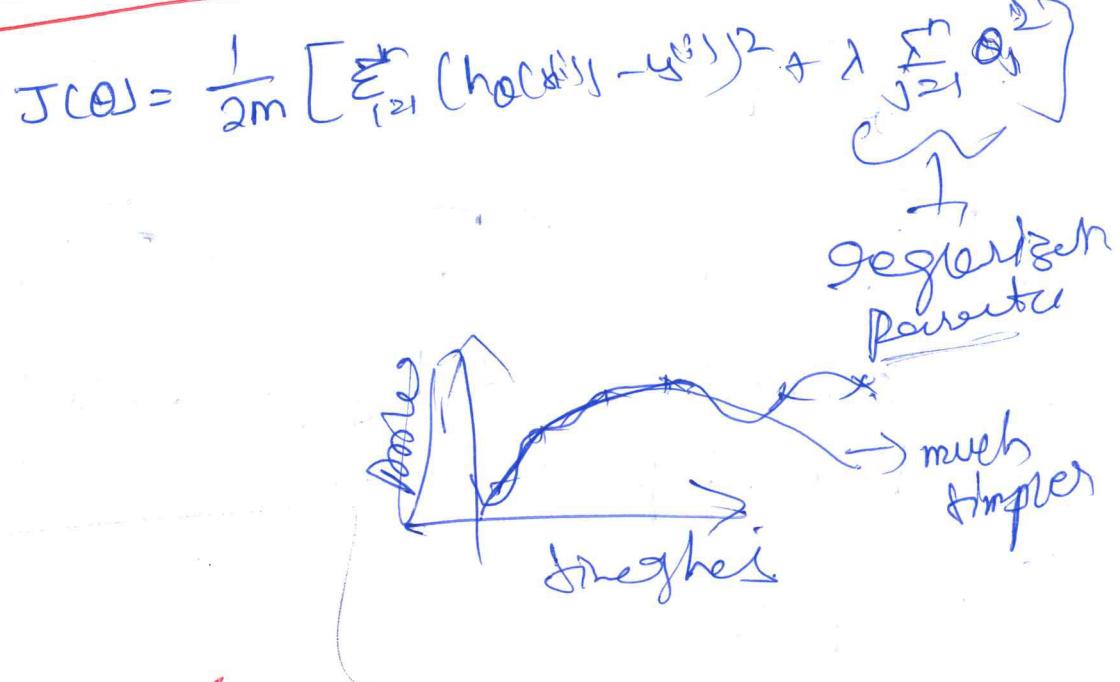
$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)})^2 + 1000\theta_2^2 + 10000\theta_4^2)$$

$$\theta_3 \approx 0, \theta_4 \approx 0$$

Small values of $\theta_0, \theta_1, \dots, \theta_n$

less overfit

Regularized



of Regularized Linear

Gradient descent
over θ 's

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

$$j = 0, 1, 2, 3, \dots, n$$

Gradient descent
linear steps

$$\theta_0 := \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})$$

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

$$= y^{(1)} - y^{(n)}$$

$$j = 1, 2, 3, \dots, n$$

regularized

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_i (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

$$\theta_j = \theta_j - \alpha \left[\frac{1}{m} \sum_i (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

Regularized J(theta)

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

$$\left(1 - \alpha \frac{\lambda}{m} \right) < 1$$

$$\theta_j = \theta_j - \alpha \cdot a \cdot a^T \theta_j$$

$$\theta_j = \theta_j - \alpha \cdot a \cdot a^T \theta_j \rightarrow \theta_j \text{ by col.}$$

Normal equation

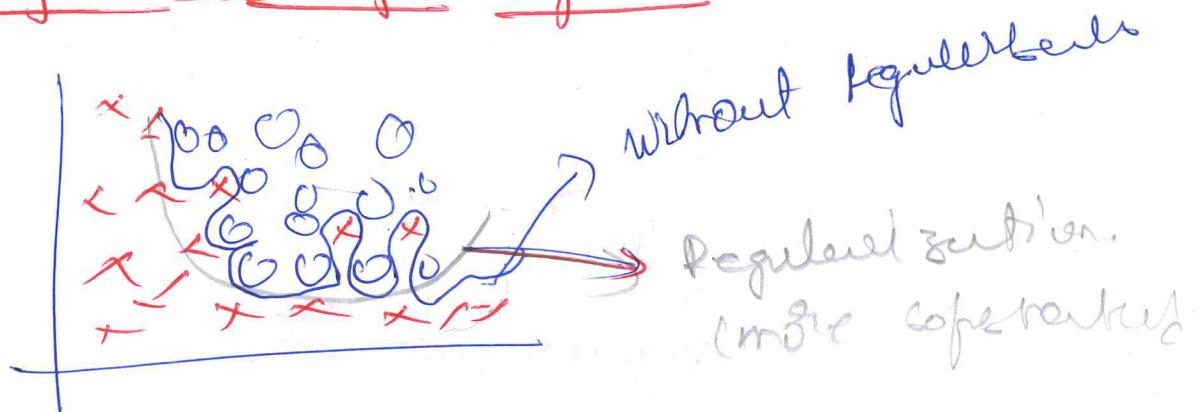
$$x = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \in \mathbb{R}^{n+1 \times m}$$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

$$\min \hat{J}(\theta)$$

$$\hat{J}(\theta) = (x^T x + \lambda I)^{-1} x^T y$$

Regularised logistic regression



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

$\theta_0, \dots, \theta_n$

How Gradient =

repeat of

$$\theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

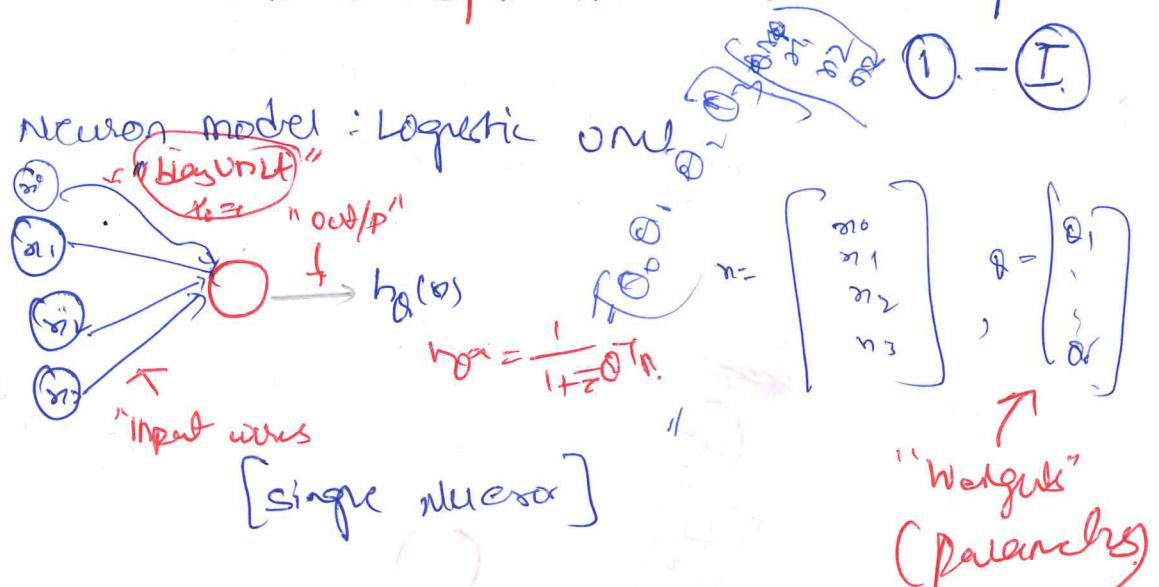
$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j$$

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

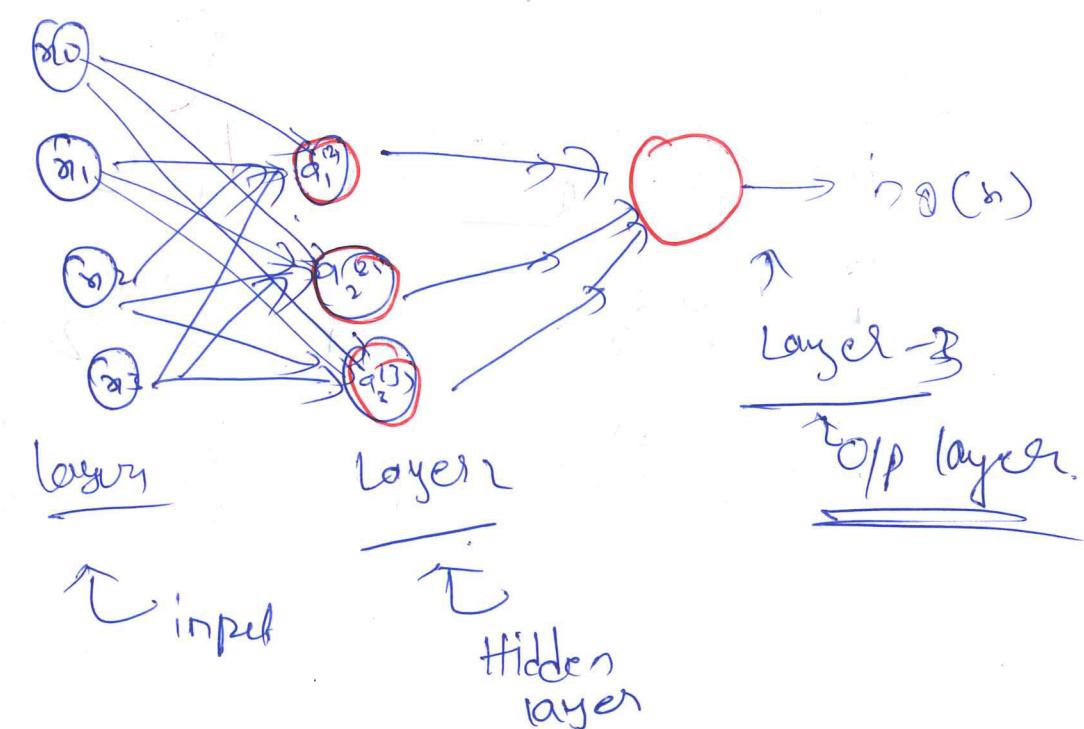
WEEK 4

NEURAL NETWORK: REPRESENTATION

— Model Representation: Model representation

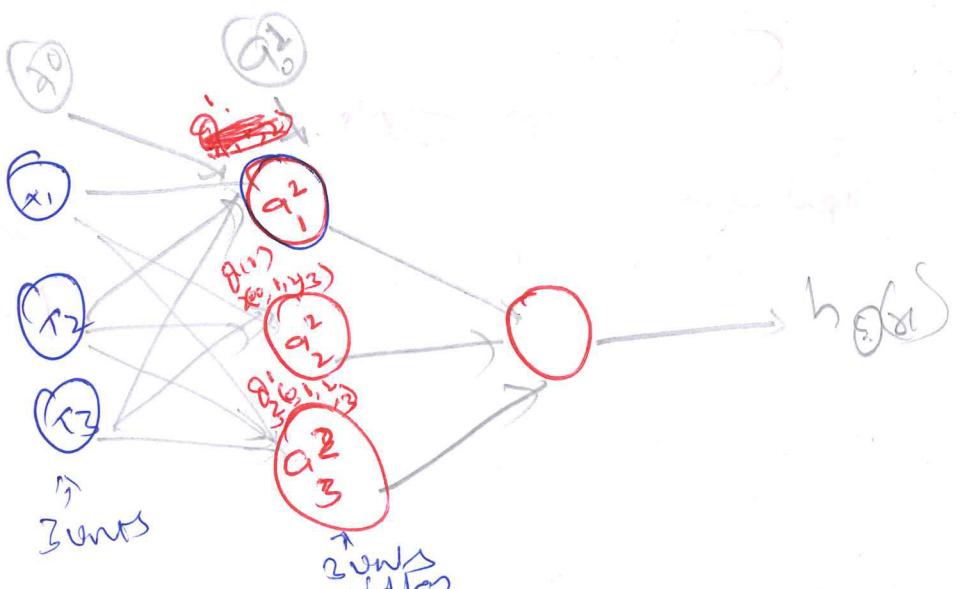


Neural Network [multiple]



a_i = "activation" of unit ; in layer j

Θ^j = matrix of weight connecting function
mapping from layer j to layer $j+1$



$$a_1^{(2)} = g(\theta_{10}^{(2)} + \theta_{11}^{(2)}x_1 + \theta_{12}^{(2)}x_2 + \theta_{13}^{(2)}x_3)$$

$$a_2^{(2)} = g(\theta_{20}^{(2)} + \theta_{21}^{(2)}x_1 + \theta_{22}^{(2)}x_2 + \theta_{23}^{(2)}x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(2)} + \theta_{31}^{(2)}x_1 + \theta_{32}^{(2)}x_2 + \theta_{33}^{(2)}x_3)$$

If in layer $j = s_j$ units

If in layer $j+1 = s_{j+1}$ units

Then $\Theta^{(j)}$ with dimension =

$$\cancel{s_j \times (s_j + 1)}$$

↓
Gvenner.

$$h_\Theta^{(j+1)} = a^{(j+1)} = g(\theta_{00}^{(j)} + \theta_{01}^{(j)}a_1 + \theta_{02}^{(j)}a_2 + \theta_{13}^{(j)}a_3)$$

Model representation \rightarrow

consider here,

$$a_1 = g(z_1^{(2)})$$

$$a_2 = g(z_2^{(2)})$$

$$a_3 = g(z_3^{(2)})$$

$$x = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}, \quad z^{(2)} = \begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{pmatrix}$$

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

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

$$\Theta^{(1)} = \Theta^{(1)} \cdot x$$

Add $a_0 = 1$

$$z^3 = \begin{pmatrix} a_0 \\ a_1 \\ a_2 \end{pmatrix}$$

$$h_{\Theta}(x) = a^3 = g(z^3)$$

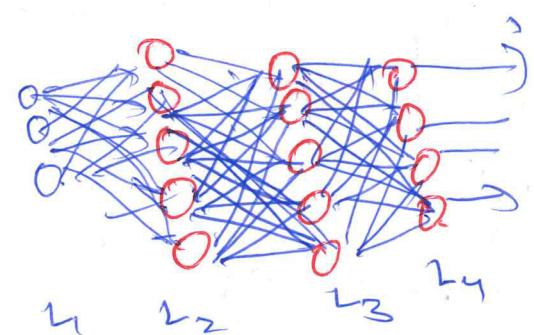
WEEK: 5

 =

LEARNING.
NEURAL NETWORKS: ~~Representation.~~

COST FUNCTION FOR FITTING Parameters

Neural Networks (Classification)



$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})\}$$

L = total no. of layers

$s_L = n$
 s_L = no. of units
(not counting bias unit)

$$s_1 = 3, s_2 = 5, s_3 = 1$$
$$s_4 = s_L = n$$

Binary classification
only one o/p
so
 $h_\theta(x) \in \mathbb{R}$
 $s_L = 1$

Multiclass

$$h_\theta(x) \in \mathbb{R}^{1C}$$
$$1C = \text{no. of o/p.}$$

$$s_L = 10, (1C \geq 3)$$

cost function from log term suggest

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

In general for a neural network we have $l^{(i)}$ o/p from $l^{(i-1)}$

$h_{\theta}(l^{(i-1)}, h_{\theta}(l^{(i)}) = l^{(i)}$

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

for 1st O/P if $i=1$, $l^{(0)}$

$$\begin{matrix} \theta^{(0)} \\ \theta^{(1)} \\ \theta^{(2)} \end{matrix} \quad y^{(0)} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

No regularization of $\theta^{(0)}$ since $\theta^{(0)}$ is bias

BACK PROPAGATION: ~~MIN~~ MIN cost function

$\min_{\theta} J(\theta)$

↳ need to compute $-J(\theta)$

$$-\frac{\partial}{\partial \theta_j} J(\theta)$$

$i =$ No. of neurons of each $x^{(i)} = x_{i,1}, x_{i,2}, \dots, x_{i,n_i}$

$j =$ No. of neurons

number $i =$ log m of each $x_j = \theta_j^{(1)}, \theta_j^{(2)}, \dots, (\theta_{j,1}, \theta_{j,2}, \dots, \theta_{j,n})$

$j =$ No. of patterns = (car, bus, toy)

Gradient computation
Given in following example (n=4)

$$\begin{matrix} \theta^{(0)} & \theta^{(1)} & \theta^{(2)} & \theta^{(3)} \\ \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

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

: Backpropagation algorithm

→ Initialize cost for error, $\delta_j^{(l)} = \text{error}$

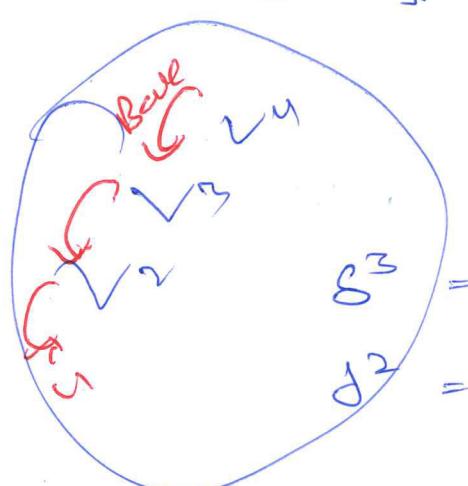
of node j in layer l .

For each output unit (layer $L=H$)

$$\delta_i^{(H)} = a_i^{(H)} - y_i$$

$$h_{\Theta}(x^{(H)} \delta_i^{(H)})$$

in vectorizing $\delta^{(H)} = a^{(H)} - y$



$$\begin{aligned}
 \delta^3 &= (\Theta^{(2)})^T \delta^{(2)} \cdot g'(z^{(2)}) \\
 \delta^2 &= (\Theta^{(1)})^T \delta^{(1)} \cdot g'(z^{(1)}) \\
 \delta^1 &= a^{(1)} - y
 \end{aligned}$$

Step of back

for one
(x, y)

=

$$\frac{\partial J(\Theta)}{\partial \Theta_{ij}} = g_j^{(l)} s_i^{(l+1)} \quad (\text{ignore } s_i \text{ if } \lambda=0)$$

For $(x^{(i)}, y^{(i)})$:- Algorithm

a row vector $((x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

Set $\Delta_{ij}^{(l)} = 0$ (for all i, j)

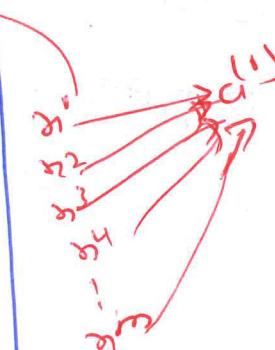
(used to compute

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$$

For $i = 1 \text{ to } m$

$$\text{Set } \delta^{(l)} = \delta^{(l)} \circ$$

Perform forward
propagation to
compute $a^{(l)}$ for
 $l = 2, 3, \dots, L$



Using $y^{(l)}$,
compute
 $\delta^{(L)} = a^{(L)} - y^{(L)}$

compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(1)}$

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \delta^{(l)} \circ \delta^{(l+1)}$$

Writing
[Agm in foest] : Backpropagation Algo

→ Training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}) \dots (x^{(t)}, y^{(t)})\}$

Set $\Delta_{ij} = 0$ (for all i, j)

For $i = 1 \dots m$

set $a^1 = x^{(i)}$

perform forward propagation for $a^{(l)}$
for $l = 2, 3, \dots L$

using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

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

→ After this

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \gamma \Theta_{ij}^{(l)} \text{ if } \gamma \neq 0$$

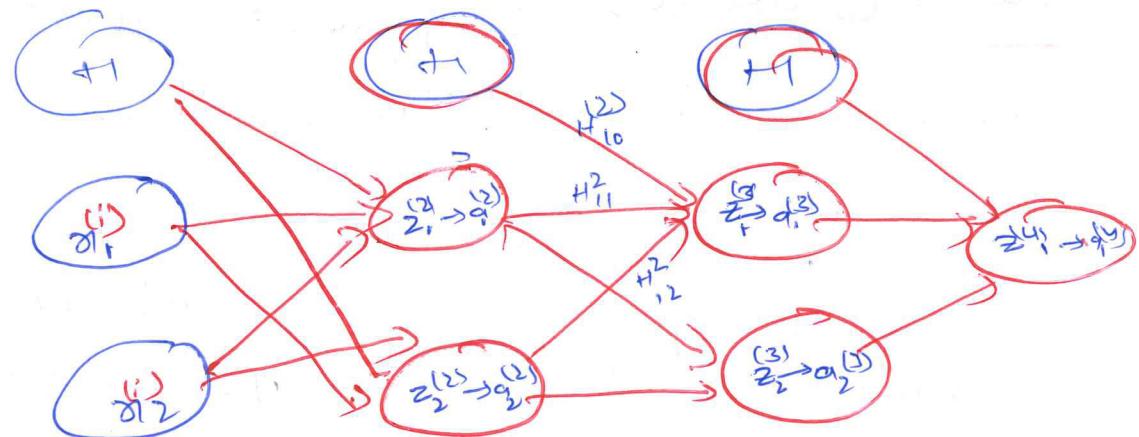
$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \left. \begin{array}{c} \frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(l)}} \\ \Delta_{ij}^{(l)} \end{array} \right\}$$

Out

~~Out of Context~~:

Forward Propagation [Input units]

Understands Back Propagation



$(x^{(i)}, y^{(i)})$

$$z_1^{(3)} = \Theta_{10}^{(2)} x_1 + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)}$$

→ In general $\delta_{ij}^{(l)} = \text{"error of cost. } a_j^{(l)} \text{ (until in layer l)}$

→ Formally- $\delta_{ij}^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$ (for $j \geq 0$),

→ $\text{cost}(i) = y^{(i)} \log h_{\Theta}(a^{(i)}) + (1-y^{(i)}) \log h_{\Theta}(a^{(i)})$ where

Back propagation In Python

Implementation note: unrolling parameters.

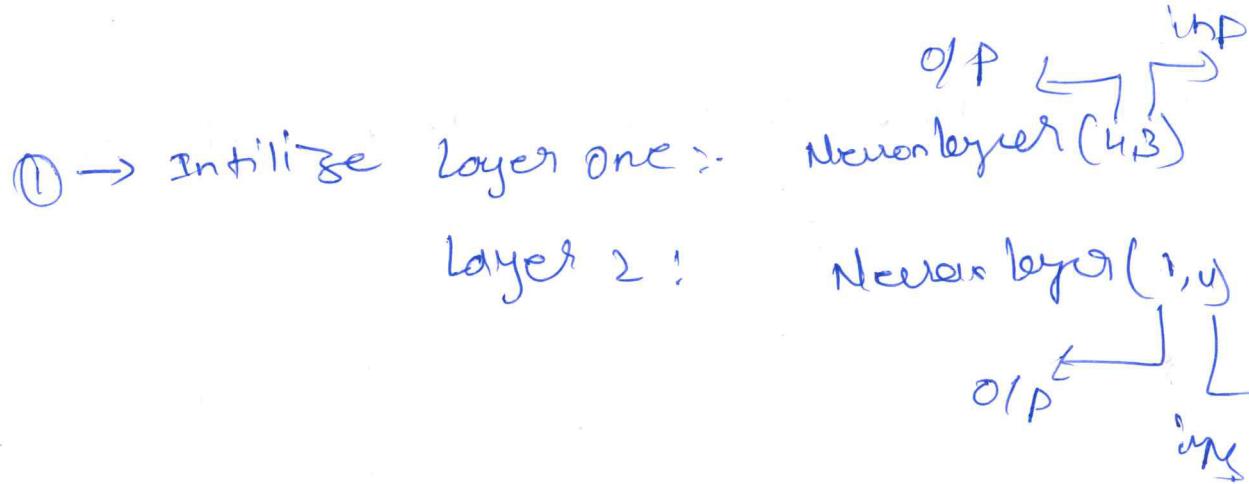
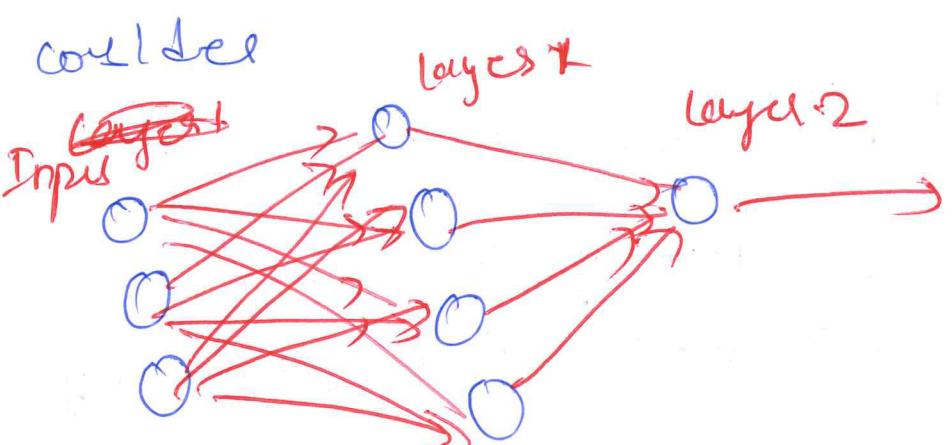
Out of Course: Implementing multilayer neural networks. In general Sing

Single layer causes error when there

is large input data because it cannot handle to me. O/P having size

so.

steps for implementing



- ① → Initialize layer one: Neuron layer (4,3)
- Layer 2: Neuron layer (3,1)
- ② → Initialize training set ^{Input} matrix
- ③ → Initialize training set output matrix
- ④ → Initialize random weights
- ⑤ → ~~Find O/P~~ → find error b/w actual L & O/P and obtain ^{Input} (during)
- B feed this ans
- ⑥ → find O/P from L₂ by feed day L₁ input to M₂

- Ⓐ → find real L_2 and $\partial/\partial L_2$
cross
- Ⓑ → Use error from $-L_2'$ time
 ΔL_2 , then ~~for~~ ΔL_1
Can be found using $\Delta \text{sign} \& \text{further}$
derivative and input.
- Ⓒ → Use ΔL_2 to find ΔL_1
- Ⓓ → Adjust weights
- Ⓔ → Again perform Ⓑ

① Debugging Learning Algorithm:

When new test is performed. We find that
There is no huge difference between
actual & predicted

We do follow:

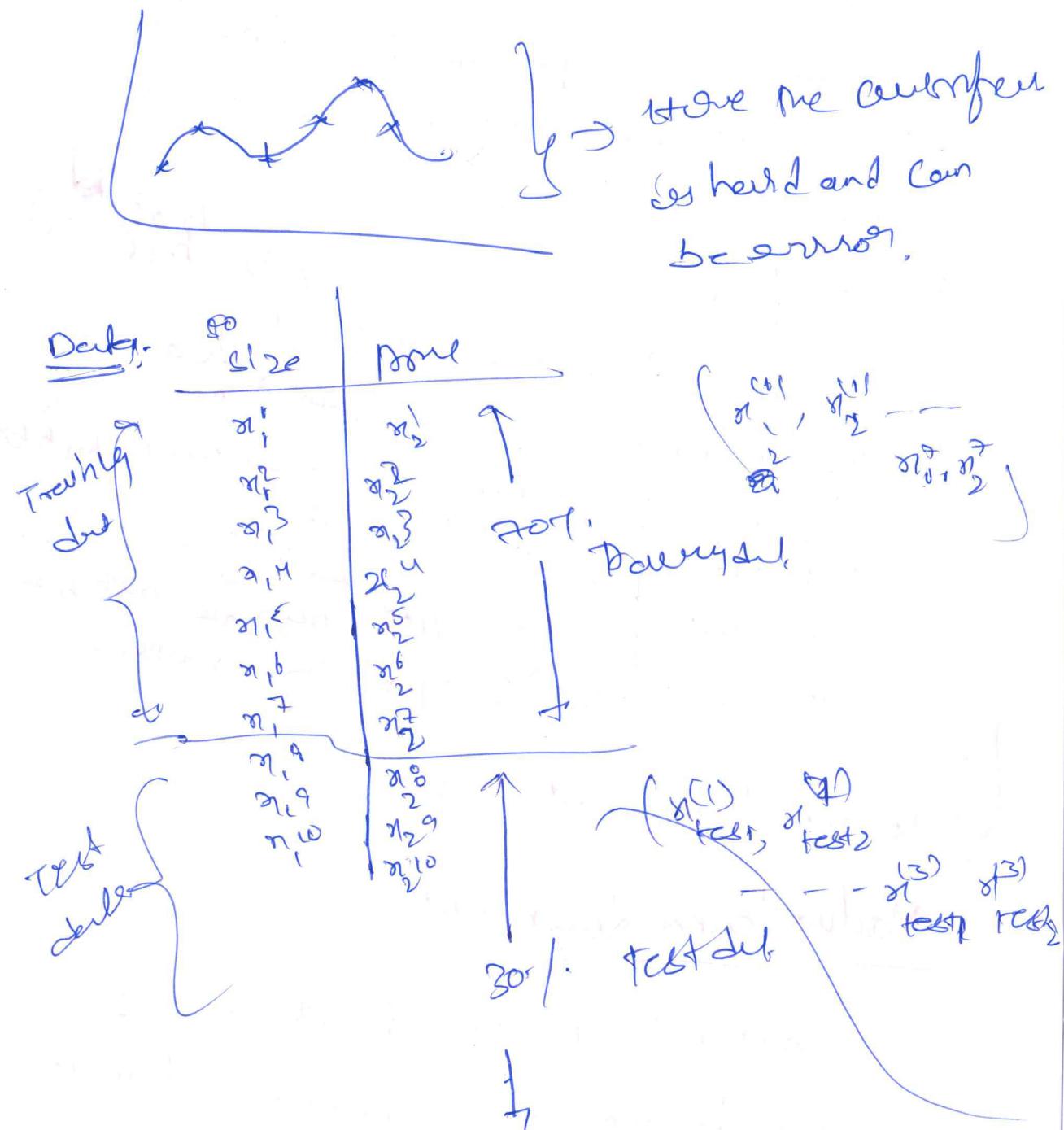
- fix more training set. ✗ (bad idea)
- Try smaller sets of features → ~~high variance~~ ^{fixed by}
- try getting additional features → ~~fixed by~~ ^{high bias}
- Try adding polynomial features ($x_1^2, x_1x_2, x_1x_3, \dots$) ^{fixes} ~~high bias~~
- Try deleting λ → ~~fixes~~ high bias
- Try increasing λ → ~~fixes~~ high variance.

↳ so we find what to do.

Machine Learning Diagnosis

Diagnostic: finding if it isn't working and gives guidance on how to improve it.

Evaluating your dataset



it gave me a wavy line
so hard and can
be error.

↑
80% training

↓ 10% cross validation

↓ 20% test

SO →

Training error:

$$J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} (h_{\theta}(x_i^{(1)}) - y_i^{(1)})^2$$

Cross validation error:

$$J_{\text{CV}}(\theta) = \frac{1}{2m_{\text{CV}}} \sum_{i=1}^{m_{\text{CV}}} (h_{\theta}(x_i^{(2)}) - y_i^{(2)})^2$$

Test error

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_i^{(3)}) - y_i^{(3)})^2$$

Model selection

$$1. h_{\theta}(x) = \theta_0 + \theta_1 n$$

$$2. h_{\theta}(x) = \theta_0 + \theta_1 n + \theta_2 n^2$$

⋮

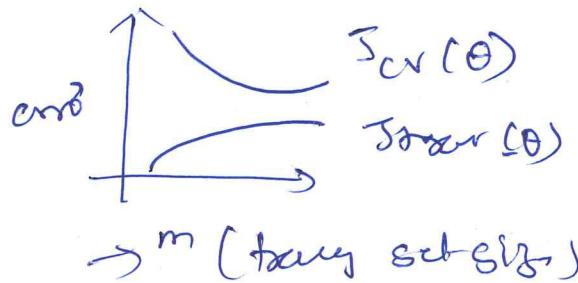
$$! h_{\theta}(x) = \theta_0 + \theta_1 n - \theta_2 n^2$$

Choose best having $\min(J_{\text{CV}})$

→ APPLY $J_{\text{test}}(\theta)$

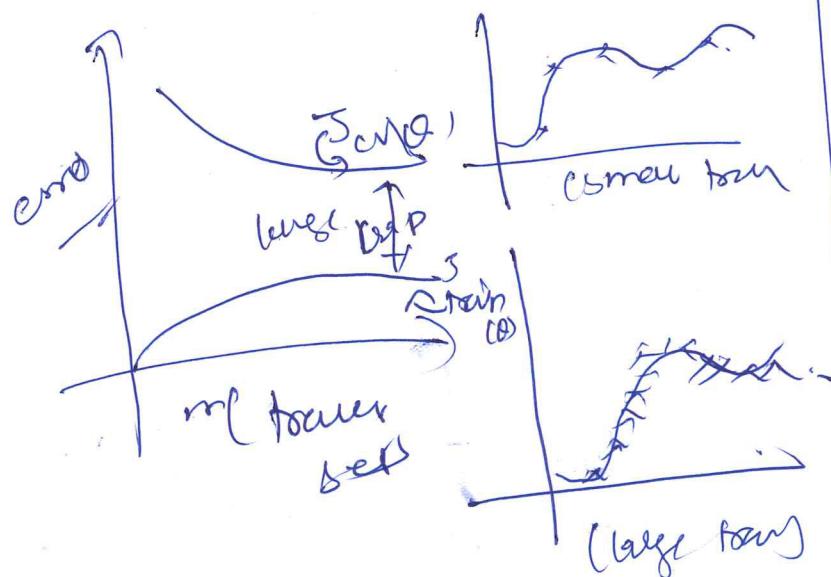
Learned levels :-

$J_{\text{train}}(\theta)$, $J_{\text{cv}}(\theta)$



For High bias

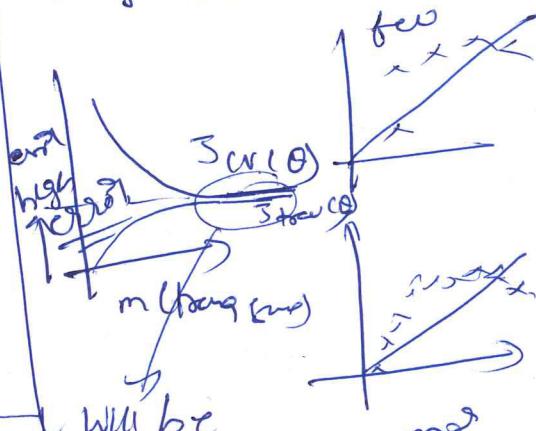
$$h_{\theta \text{ (BS)}} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots - \theta_{100}^{100}$$



For High variance

$$\text{Supr. bias} = \theta_0 + \theta_1 x_1$$

Requires lot more training



Small Network
Network

computationally
cheaper

"Large network
networks"

computationally
more expensive

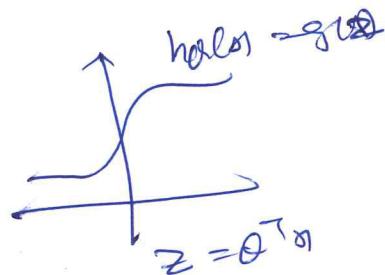
use regularization (λ)
to address overfitting,

SUPPORT VECTOR MACHINES

I) OPTIMIZATION objective :-

Alternate view of logistic regression

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

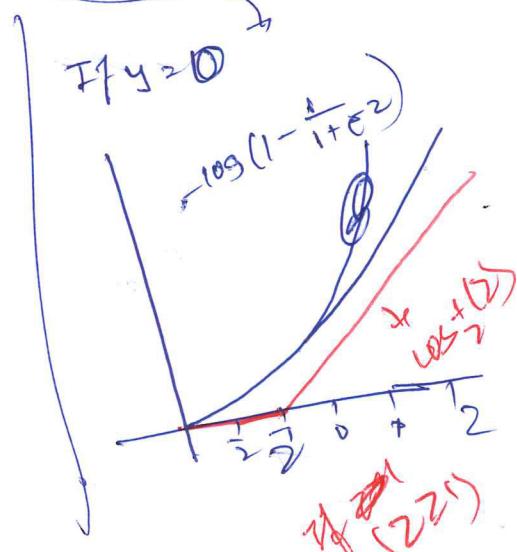
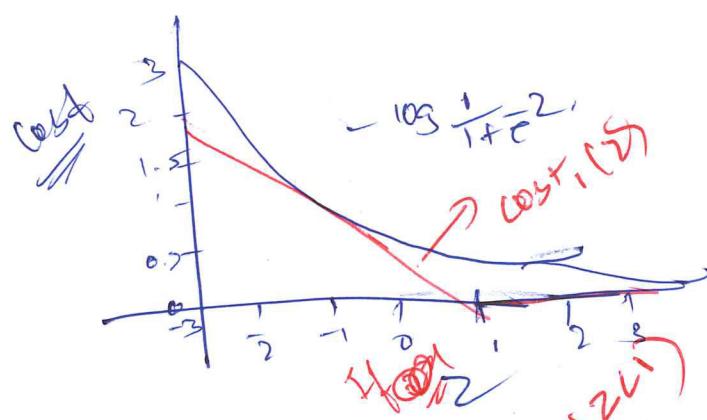


If $y=1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x \gg 0$
 $y=0$.. $h \approx 0$, $\theta^T x \ll 0$

* cost function of logistic:

$$\begin{aligned} &= -y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x)) \\ &= -y \log \frac{1}{1+e^{\theta^T x}} + (1-y) \log \left(1 - \frac{1}{1+e^{\theta^T x}}\right) \end{aligned}$$

If $y=1$, ($\text{want } \theta^T x \gg 0$)



SVM:

logistic Regression:

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

Support Vector Machine:

$$\begin{aligned} \min_{\theta} & \frac{1}{m} \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] \\ & + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \end{aligned}$$

SVM Hypothesis:

$$\min_{\theta} c \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

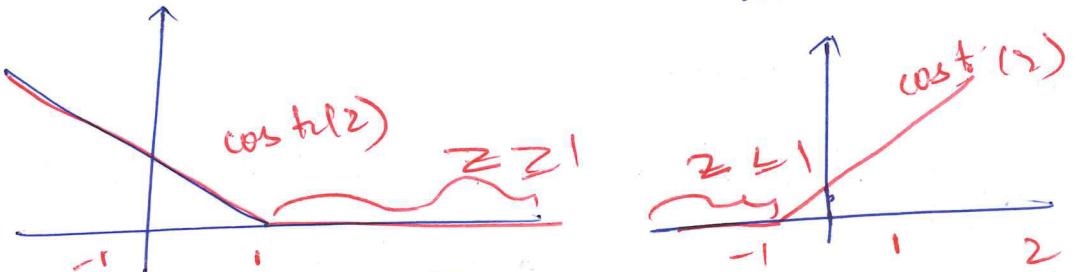
Hypothesis:

$$h_{\theta}(x) = \begin{cases} 1, & \text{if } \theta^T x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

SG

SVM: Large Margin Classification

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



If $y=1$, we use $\theta^T x_i$ (not just ≥ 0), $\theta^T x_i \neq 0$
 If $y=0$, we use $\theta^T x_i - 1$ (not just ≤ 0) $\theta^T x_i \neq -1$

if $y \neq 0$ does rule

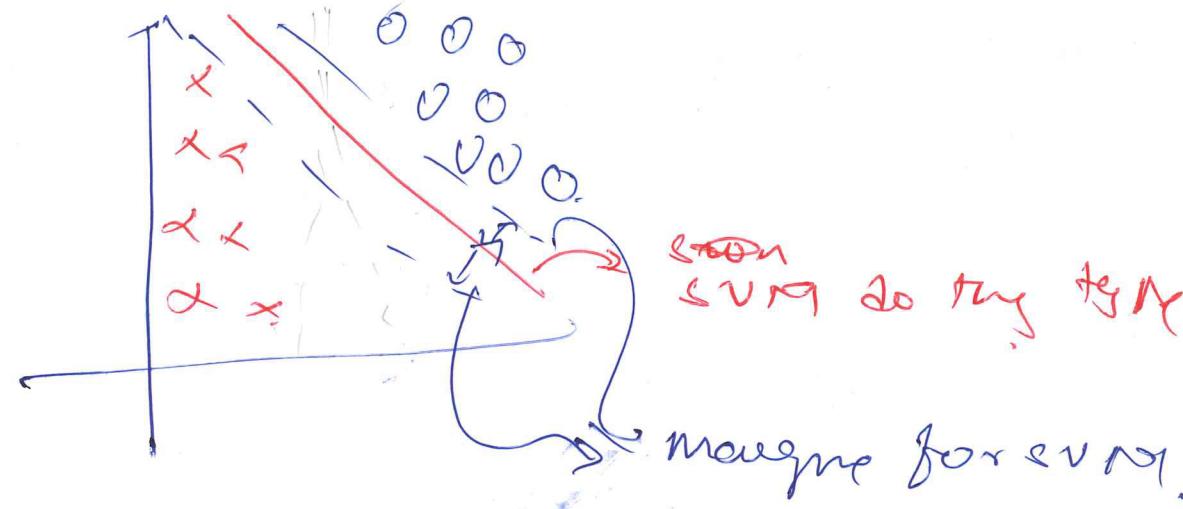
SVM :- Boundaries

$$\theta^T x = 0$$

$$\min C D + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$H \quad \text{If } y^{(i)} = 1 \\ \cancel{\theta^T x_i + \theta_0} + \theta^T x_i \geq 1 \\ \theta^T x_i \geq 1$$

If $y^{(i)} = 0$
 $\theta^T x_i \leq 0$



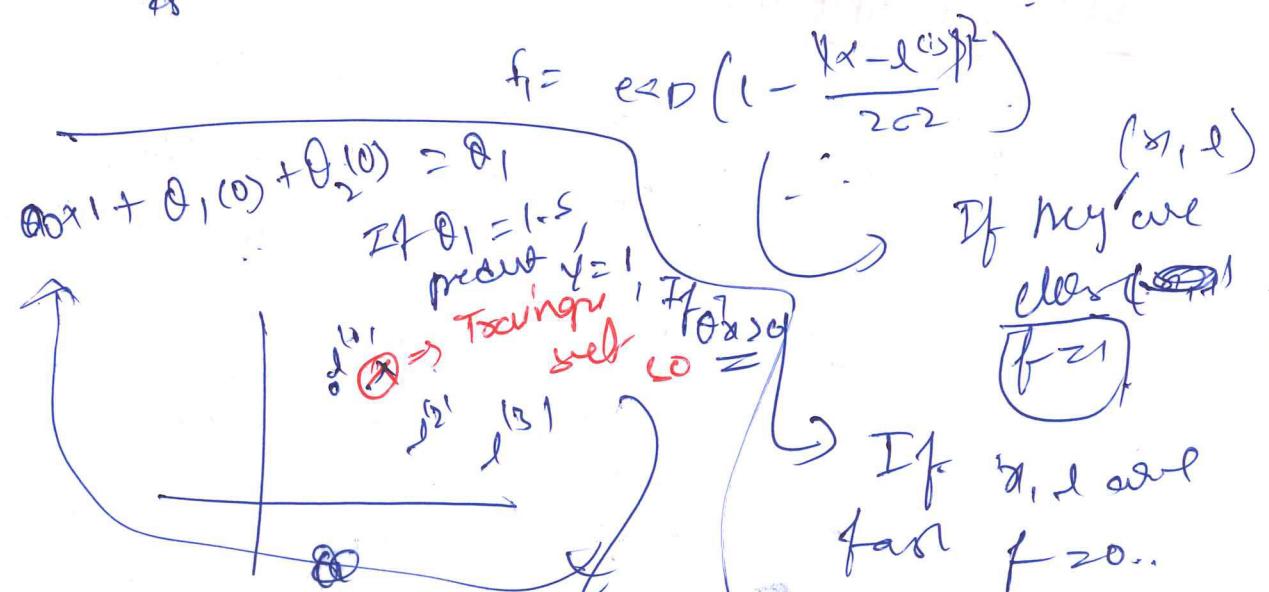
SVM with Kernels

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

→ predict " $y=1$ " if $\theta^T f \geq 0$

$$\theta_0 + \theta_1 f_1 + \dots + \theta_m f_m$$

$$\begin{aligned} f_0 &= x_0 \\ f_1 &= x_0 x_1 \\ f_2 &= x_0 x_2 \\ &\vdots \end{aligned}$$



$$\text{Training} \underset{\theta}{\min} C \sum_{i=1}^m g^{(i)} \text{cost}_i(\theta^T \phi^{(i)}) + (\frac{1}{2}) \sum_{j=1}^{n+m} \theta_j^2$$

$$g_j(\theta_j) = \theta^T \phi_j$$

$$\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}, \text{ (ignore } \theta_0)$$

$$\Rightarrow \theta = \theta^T m \theta, \left\{ \text{making me equal for word.} \right.$$

To run SVM effectively

SVM Parameters:

$C (= \frac{1}{\lambda})$, Large C = lower bias, high variance
Small C = higher bias, low variance

σ^2 Large σ^2 : forces f: very smooth, higher bias, low variance.

Logistic Regression vs SVM

n = number of features, ($\propto E2^{n+1}$), m = number of training examples

If n is large (relative to m) :

→ Use logistic regression & SVM without kernel (linear kernel)

→ If n is small, m is intermediate:

$$\rightarrow \text{use SVM with Gaussian Kernel} \\ \therefore f = -\exp\left(-\frac{(x-\mu)^2}{\sigma^2}\right)$$

→ If n is small, m is large

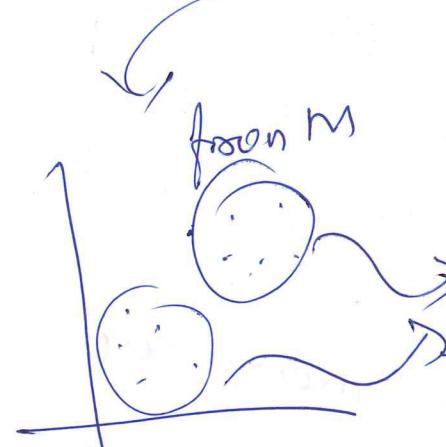
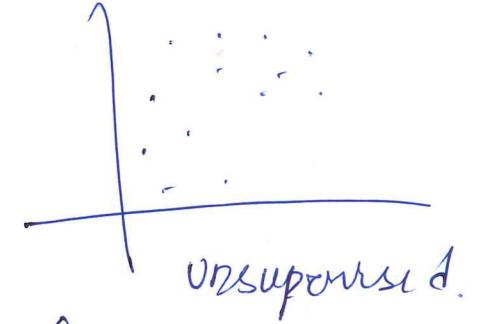
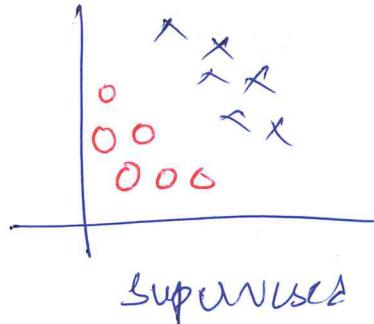
→ overfit/ add more features, then use logistic regression & SVM without a kernel.

→ Newton's method works for most of these settings but may be slower to train.

Unsupervised Learning

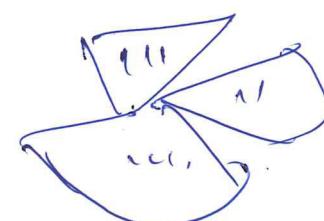
Unsupervised Learning (No Labeling as in supervised)

Introduction to-

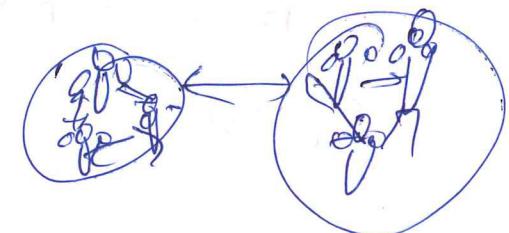


from m

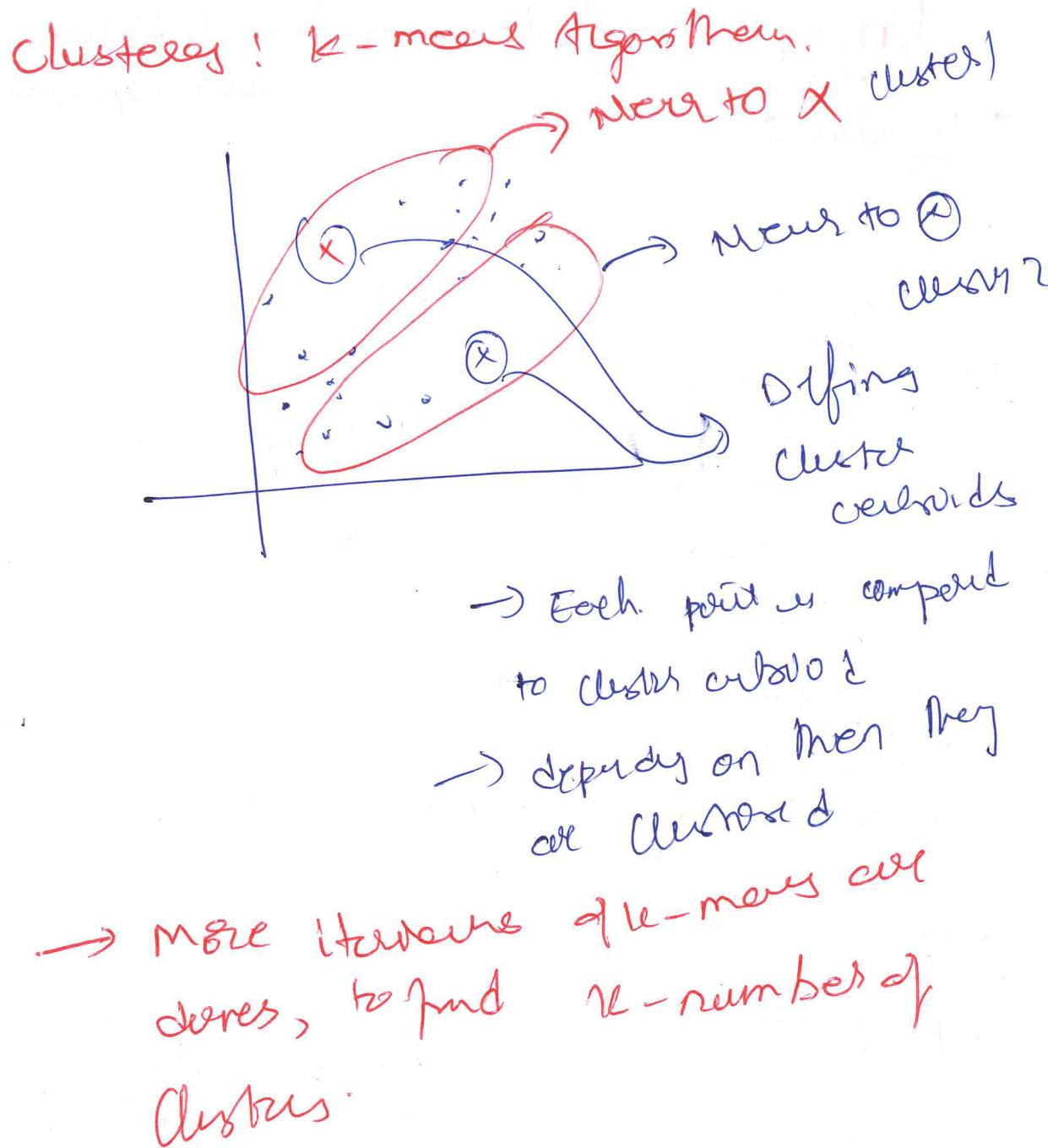
cluster algorithm.



App of clusters



Neural network.



K-means algorithm.

Input:

- k (number of clusters)
- Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

$x^{(i)} \in \mathbb{R}^n$ ($d_{\text{obs}} = 1$, dimension)

K-means algorithm

Random initialize K clusters $u_1, u_2, u_3, \dots, u_k$

~~repeat loop for k~~

repeat {

for $i = 1$ to m

$c^{(i)} = \text{index (from 1 to } K) \text{ of cluster centered close to } x^{(i)}$

for $k = 1$ to K

$u_k = \text{average/mean of points assigned to cluster } k$

~~choose
move
centroid~~

~~min J w.r.t.
 u_1, \dots, u_K~~

~~repeat~~

~~until converge~~

~~or until ID. update~~

~~or update
centroid~~

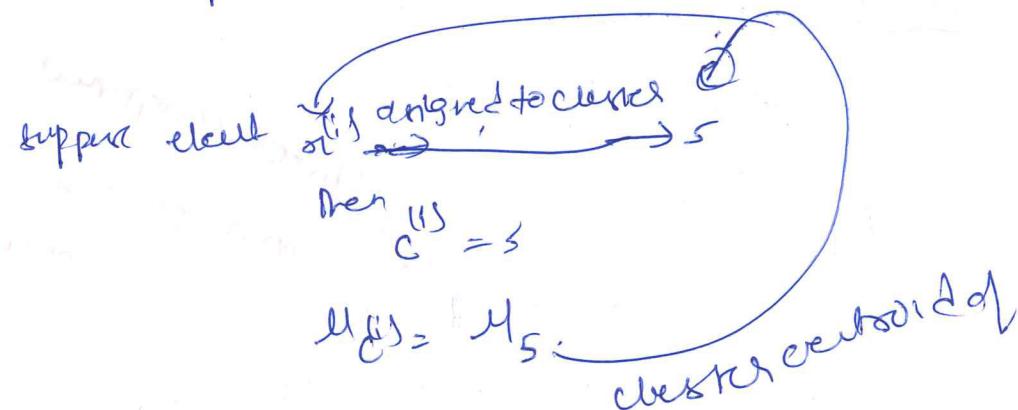
K-Means optimization Objective

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

$x^{(i)}$ is currently assigned.

μ_k = cluster centroid $\in \{\mu_1, \dots, \mu_K\}$ $k \in \{1, 2, \dots, K\}$

$\mu_c^{(i)}$ = cluster centroid of cluster to which example
example $x^{(i)}$ has been assigned



Optimizierter objektivus

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

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

distance between

centroid and object
has to be minimize & next

they can be almost clustered same.

K-Means Initialization

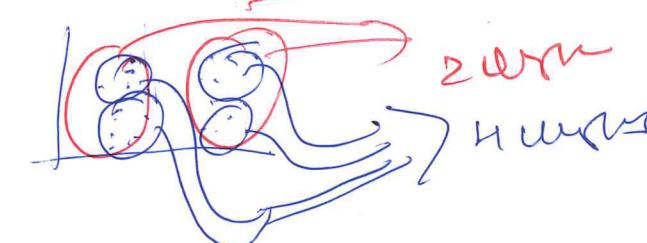
~~see driver for notes~~

random initialization

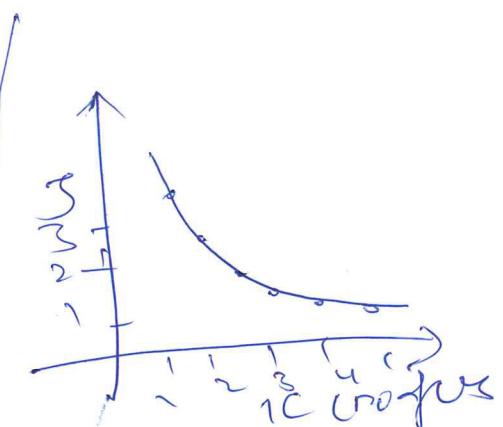
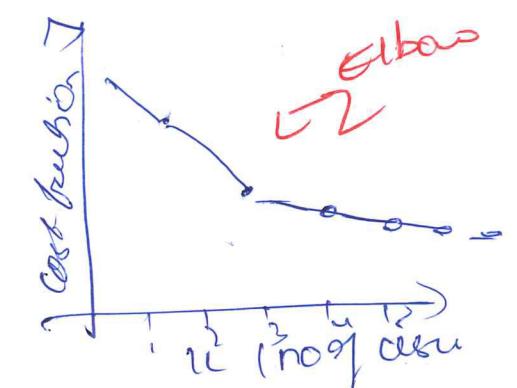
- should $n \geq m$
- randomly pick ~~one~~ m training examples

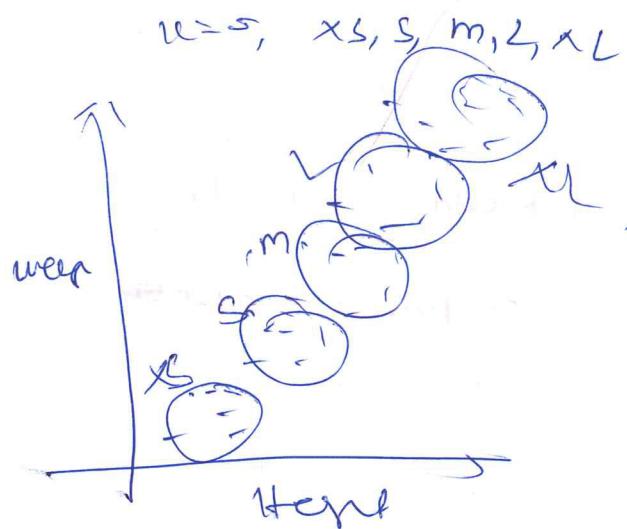
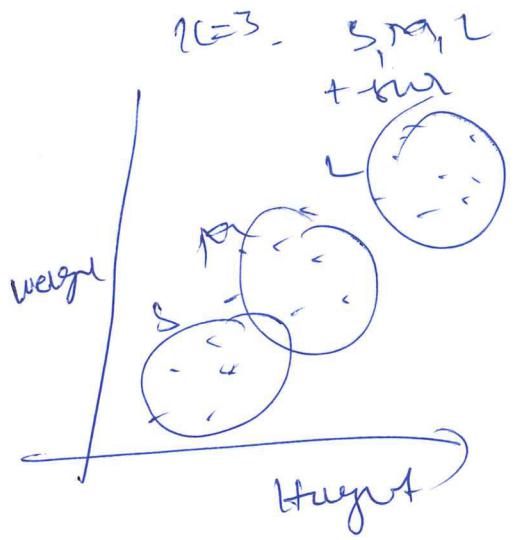
- set μ_1, \dots, μ_m equal to
these n samples

choose m number of clusters

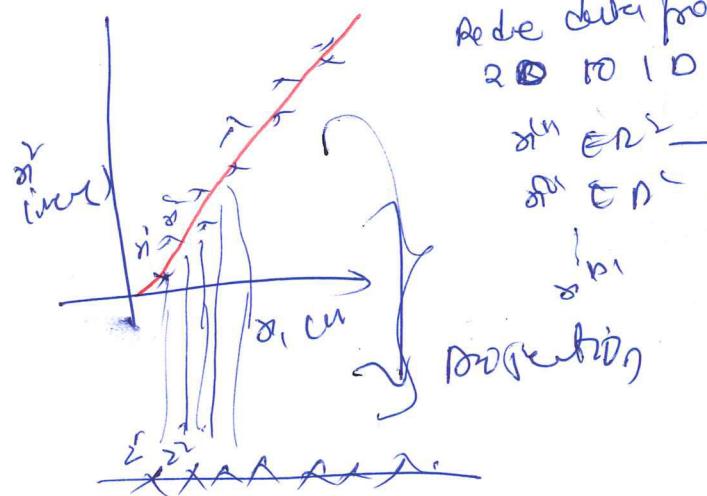


Elbow method





Data compression

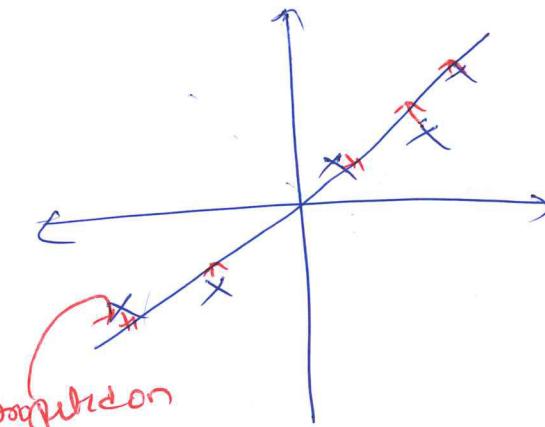


Reduce data from
2D to 1D

$$\begin{matrix} \text{in } \mathbb{R}^n & \rightarrow & \text{in } \mathbb{R}^1 \\ \text{in } \mathbb{C}^n & \rightarrow & \text{in } \mathbb{C}^1 \\ \text{in } \mathbb{R}^m & \rightarrow & \text{in } \mathbb{R}^1 \end{matrix}$$
 $\rightarrow \text{PCA}$

Dimensionality Reduction: Principal Component Analysis.

PCA: Formulation



PCA - Reduce 2D - to - 1D with min. projection.

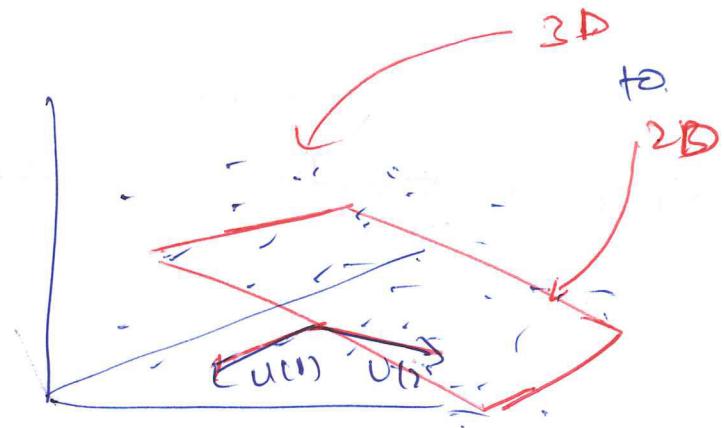
Find direction (a vector $u \in \mathbb{R}^n$).

→ General Case

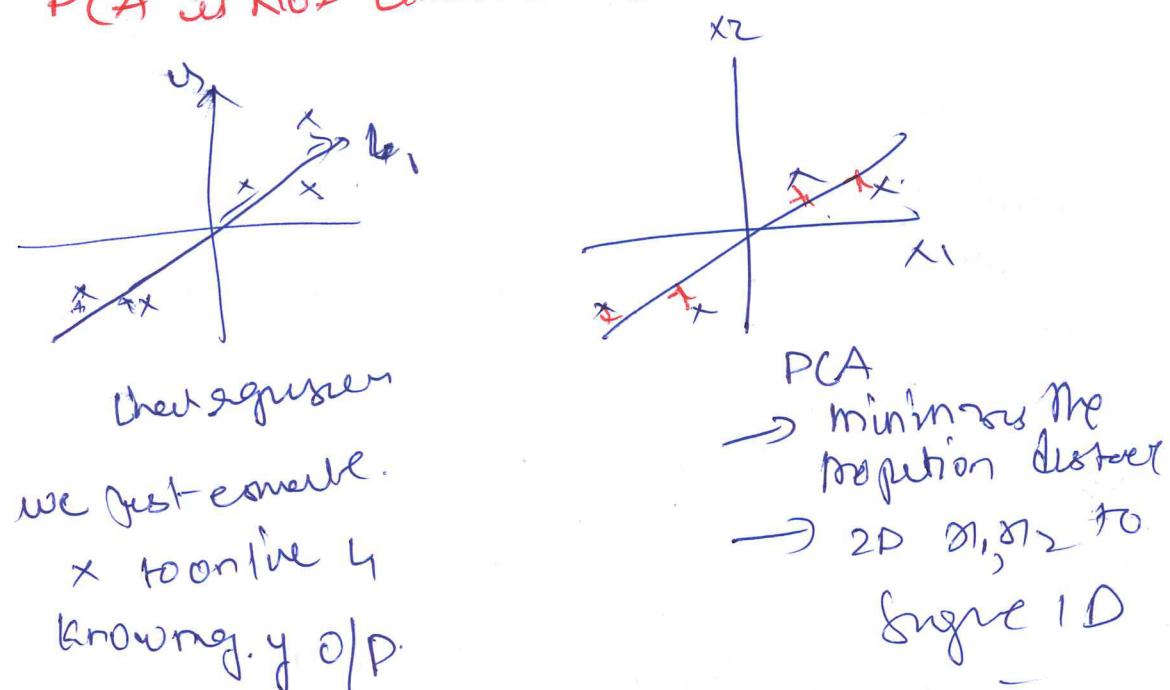
Reduce n -dimension to k -dimensions: find k -vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

$x \in \mathbb{R}^n$
 \rightarrow low dimension
 \rightarrow PCA = does the
minimum of sum
of squares of
projected distance

20



→ PCA is NOT linear regression



PCA Algorithm

Data Preparation of Data set: $\{x^{(1)}, \dots, x^{(m)}\}$
 $\bar{x}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$ scaling term (mean Normalization)

Replace each $x_j^{(i)}$ with $x_j^{(i)} - \bar{x}_j$

- = If different features on different scale (e.g. x_1 , size of house, x_2 = number of bedrooms), scale features to have comparable magnitude

* → Reduce dimensions from n-D to 1-D
 compute 1st covariance matrix

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$$

→ Compute eigen vectors $\{\cdot\}$:

$$[U, S, V] = \text{SVD}(\Sigma)$$



Singular
Value
Decomposition

$$U = \begin{bmatrix} | & | & | \\ U^{(1)} & U^{(2)} & \dots & U^{(m)} \\ | & | & | \end{bmatrix}$$

$U \cdot x \in \mathbb{R}^n \rightarrow \mathbb{R}^k$

~~take~~
 $x \in \mathbb{R}^n \Rightarrow z \in \mathbb{R}^{k'}$



Mathematical Proof Beyond the Scope
on = n x k
 U^T_{order}

$(U^T_{\text{order}})^T x^* = \begin{bmatrix} N_{11} & \dots & N_{1k'} \\ \vdots & \ddots & \vdots \\ N_{k'n} & \dots & N_{kk'} \end{bmatrix} \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(k')} \end{bmatrix}$

Summary of PCA

Algorithms

$$\Sigma = \frac{1}{m} \sum (x^{(i)}) (x^{(i)})^T$$

$$[U, S, V] = \text{svd}(\Sigma)$$

$$U^{\text{order}} = U(:, 1:k')$$

$$Z = U^{\text{order}} * x;$$

$$X = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(m)} \end{bmatrix}$$

$$\hat{x}_{\text{app}} = \frac{1}{m} \times X$$

Reconstruction from compressed representation!

$$Z = U^T_{\text{order}} * x$$

$$x = \underset{\text{approx}}{\text{approx}} U^T_{\text{order}} * Z$$

on = n x k

choosing k' (NodePCA)

$$\text{Avg Error} : \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x^{(i)}_{\text{approx}}\|^2$$

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

* Choose k' to be Smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x^{(i)}_{\text{approx}}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01$$

→ 99% variation ~~greater~~ retained

Algorithm :- (L2 Checks)

Try PCA with $k=1$

Compute $V_{approx}, z_1^{(1)}, z_2^{(1)}, \dots, z_m^{(1)}$

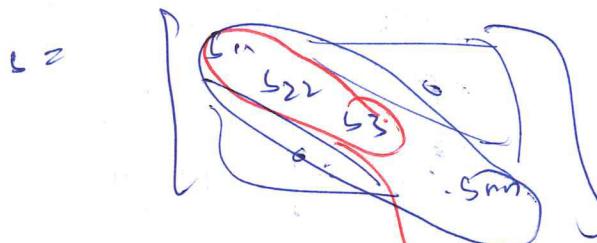
$x_i^{(1)} appr \dots x_i^{(m)}$

Check if

$$\frac{1}{m} \sum_{i=1}^m \| (x_i^{(1)} - x_i^{(1)} appr) \|^2 \leq 0.01$$

$$\frac{1}{m} \sum_{i=1}^m \| x_i^{(1)} \|^2$$

$[V, S, N] = \text{WV (kigns)}$



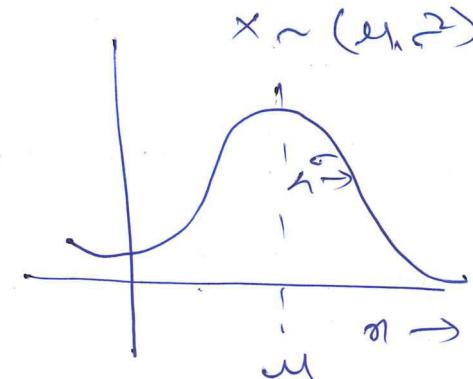
for que 1c

$$1 - \frac{\sum_{i=1}^m s_{ii}}{\sum_{i=1}^m s_{ii}^2} \leq 0.01$$

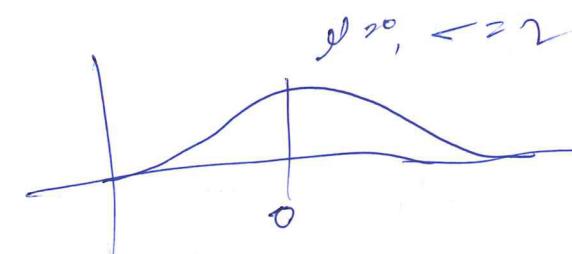
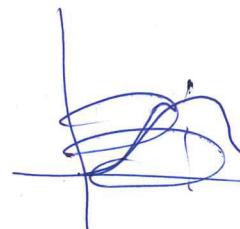
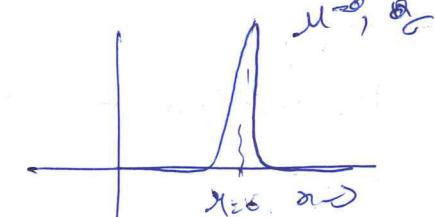
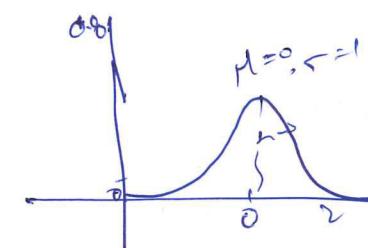
Density Estimation :

Gaussian distribution

if $x \sim \mathcal{N}(\mu, \sigma^2)$, σ is distribution with mean μ , variance σ^2



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



Training set : $\{x_1^{(1)}, \dots, x_n^{(1)}\}$

Each sample $x_i \in \mathbb{R}^d$

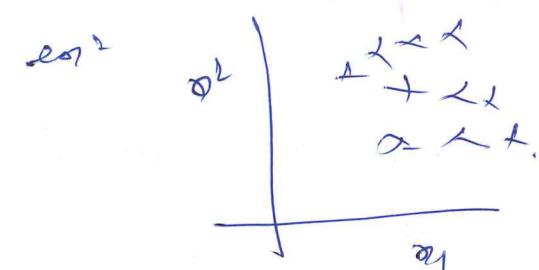
(x_i, u_i, v_i)

$$P(\sigma) = P(x_1, u_1, v_1) \cdot P(x_2, u_2, v_2) \cdots P(x_n, u_n, v_n)$$

Every of them has different
Gaussian parameters

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

Anomaly of P(D)LE



Anomaly detection algorithm (Assume Gaussian distribution as normal fault is out of left.)

1. Choose feature x_i , that is general responsible for fault.

2. Fit parameters $\mu_1, \mu_2, \dots, \mu_n, \sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$

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

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

3. Given new sample x , compute $P(x|\sigma)$

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

PREDICTIVE MOVIE RATINGS

Recommendation System — Recommended products based on previous searches for different types of filters for a particular user.

n_u = no of users

Problem Formulation :-

$\rightarrow r(i,j) = 1$, if user i has rated movie j
 $(0, 0 \text{ else})$

$\rightarrow y_{(i,j)} = \text{rating by user } j \text{ for movie } i$
 (if defined : if $r(i,j) = 1$)

→ $\theta^{(j)}$ = parameter vector for user j

$\rightarrow \mathbf{x}^{(1)}$ = feature vector for movie!

→ For user j , movie i , predicted rating: $\hat{r}_{j,i} = \theta_0^j + \theta_1^j S_{j,i} + \theta_2^j T_{j,i}$

$m^{(j)}$ = no. of movies rated by User j

To learn $\theta^{(j)}$:

$$\min_{\theta(i)} \frac{1}{2m} \sum_{i=1}^m \sum_{j=1}^{n_i} \left((\theta(i))^T x^{(i,j)} - y^{(i,j)} \right)^2$$

↓
Constant + $\frac{\lambda}{2m} \sum_{k=1}^n (\theta_k)^2$

$y_i = b_i + \epsilon_i$

	x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
Dom 1	1	0	0	0	0	0	0	0	0	0	0
Dom 2	0	1	0	0	0	0	0	0	0	0	0
Dom 3	0	0	1	0	0	0	0	0	0	0	0
Dom 4	0	0	0	1	0	0	0	0	0	0	0
Dom 5	0	0	0	0	1	0	0	0	0	0	0
Dom 6	0	0	0	0	0	1	0	0	0	0	0
Dom 7	0	0	0	0	0	0	1	0	0	0	0
Dom 8	0	0	0	0	0	0	0	1	0	0	0
Dom 9	0	0	0	0	0	0	0	0	1	0	0
Dom 10	0	0	0	0	0	0	0	0	0	1	0

Optimization problems
Rebalancing movie selection

Optimization problems

To learn $\Theta^{(j)}$ (Parameters for v_{0j}, v_{1j}):

$$\min_{\theta(i)} \frac{1}{2} \sum_{j: \delta(i,j)=1} \left((\theta(i))^T x(j) - y(i,j) \right)^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta(k))^2$$

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_w)}$:

for convexity

$$\min_{\theta_1, \dots, \theta_n} \frac{1}{2} \sum_{j=1}^n \sum_{i=(i,j) \geq 1} \left((\theta^{(ij)}_j)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{l=1}^n \left(\theta^{(lj)}_j \right)^2$$

so finally:

optimization algorithm:

$$\min_{\theta^{(0)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{j=1}^n \sum_{i: r(i,j)=1} ((\theta^{(0)})^T \phi^{(i,j)} - u^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{i: r(i,j)=1} (\theta_j^{(0)})^2$$

$\mathcal{J}(\theta_1, \dots, \theta_n)$

gradient descent update:

$$\theta_{1c}^{(i)} := \theta_{1c}^{(i)} - \alpha \sum_{j: r(i,j)=1} ((\theta^{(0)})^T \phi^{(i,j)} - u^{(i,j)}) \phi_{1c}^{(i,j)} \quad \text{for } u=0$$

$$\theta_{1c}^{(i)} := \theta_{1c}^{(i)} - \alpha \sum_{j: r(i,j)=1} ((\theta^{(0)})^T \phi^{(i,j)} - u^{(i,j)}) \phi_{1c}^{(i,j)} + \lambda \theta_{1c}^{(i)} \quad \text{(for } i \neq 0)$$

$$[00] \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \\ \theta^{(2)} \end{bmatrix} = \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \\ \theta^{(2)} \end{bmatrix}$$

$$[00] \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \end{bmatrix} = 0, \quad [03] \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \end{bmatrix} = 0.5$$

$$[05] \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \end{bmatrix} = 2.5$$

$$\theta_1 = 1.5 \\ \theta_1 = \frac{0.5 + 2.5}{3}$$

Collaborative Filtering

Given $\theta^{(0)}, \dots, \theta^{(n)}$ (and movie ratings)
can estimate $\theta_1, \theta_2, \dots, \theta^{(n)}$

Given $\theta^{(0)}, \dots, \theta^{(n)}$

can estimate $\theta^{(0)}, \dots, \theta^{(n)}$

Graph $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow \dots$

Collaborative Filtering optimization objective

Given $\theta^{(1)}, \dots, \theta^{(n)}$, estimate $\theta^{(1)}, \dots, \theta^{(n)}$

$$\min_{\theta^{(1)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{j=1}^n \sum_{i: r(i,j)=1} ((\theta^{(i)})^T \theta^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta^{(j)}_k)^2$$

Given $\theta^{(1)}, \dots, \theta^{(n)}$, estimate $\theta^{(1)}, \dots, \theta^{(n)}$:

$$\min_{\theta^{(1)}, \dots, \theta^{(n)}} \frac{1}{2} \sum_{i=1}^m \sum_{j: r(i,j)=1} ((\theta^{(i)})^T \theta^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (\theta^{(i)}_k)^2$$

→ Minimize both simultaneously.

$$J(\theta^{(1)}, \dots, \theta^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$$

$$= \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(i)})^T \theta^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (\theta^{(i)}_k)^2 + \frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\theta^{(j)}_k)^2$$

$$\min_{\theta^{(1)}, \dots, \theta^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}} J(\theta^{(1)}, \dots, \theta^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$$

Collaborative Filter algorithm

1. Initialize $\theta^{(1)}, \dots, \theta^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}$ to small random values.
2. Minimize $J(\theta^{(1)}, \dots, \theta^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$ Use gradient descent (or any advanced optimizers algorithm) For Eg. for

every $j=1, \dots, n$ $i=1, \dots, m$

$$\theta^{(j)}_k := \theta^{(j)}_k - \alpha \left(\sum_{i: r(i,j)=1} ((\theta^{(i)})^T \theta^{(j)} - y^{(i,j)}) \theta^{(i)}_k + \lambda \theta^{(j)}_k \right)$$

$$\theta^{(i)}_k := \theta^{(i)}_k - \alpha \left(\sum_{j: r(i,j)=1} ((\theta^{(i)})^T \theta^{(j)} - y^{(i,j)}) \theta^{(j)}_k + \lambda \theta^{(i)}_k \right)$$

$$\frac{d J}{d \theta^{(j)}_k}$$

Low Rank Matrix Factorization

$$\gamma = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 8 \\ ? & 4 & 0 & 8 \\ 0 & 0 & 0 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

Predicted rating

$$(\theta^{(1)})^T \theta^{(1)} + (\theta^{(2)})^T \theta^{(2)} + \dots + (\theta^{(n)})^T \theta^{(n)}$$

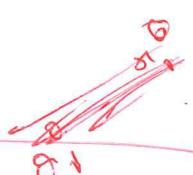
$$(\theta^{(1)})^T \theta^{(1)} + (\theta^{(2)})^T \theta^{(2)} + \dots + (\theta^{(n)})^T \theta^{(n)}$$

$$\theta^{(1)} \gamma \theta^{(1)}$$

$$\theta^{(2)} \gamma \theta^{(2)}$$

$$\vdots$$

$$\theta^{(n)} \gamma \theta^{(n)}$$



$$x = \begin{pmatrix} - & \frac{x^{(1)T}}{(x^{(1)})^T} \\ - & \frac{x^{(2)T}}{(x^{(2)})^T} \\ - & \vdots \\ - & \frac{x^{(m)T}}{(x^{(m)})^T} \end{pmatrix}$$

$$\Theta = \begin{pmatrix} - & \theta^{(0)T} \\ - & \theta^{(1)T} \\ - & \vdots \\ - & \theta^{(n)T} \end{pmatrix}$$

SUMMARY OF COURSE

- Andrew Ng

Main topics

$$\Theta^{(1)T} x^{(1)} \cdots \Theta^{(n)T} x^{(n)} = X \Theta^T$$

Low Rank Matrix factorization

ML Topics

① Supervised Learning

- Linear regression, logistic regression, Neural Networks, SVMs

② Unsupervised Learning

- K-means, PCA, Anomaly detection.

③ Special applications / special topics

- Recommendation system, large scale machine learning.

④ Advice on building a machine learning system

Bias/variance, regularization; dealing with overfitting

Algorithm selection; what to work on next; evolution of learning algorithms, learning curves, error analysis, testing analysis.

ML Applying tools

172
Court
Complex