

Machine Learning - Andrew Ng.

ML

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

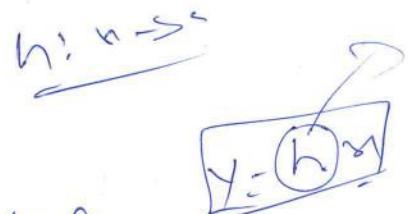
x : set of inputs of size n

$$x = (x^n)$$

y : set of outputs

$$y = y^n$$

$\rightarrow h: \mathbb{R}^n \rightarrow \mathbb{R}$, when ~~$n > 1$~~ , $y \in \mathbb{R}$



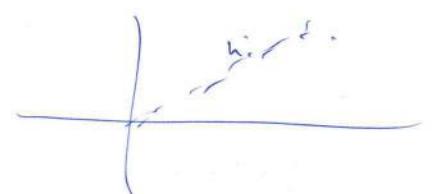
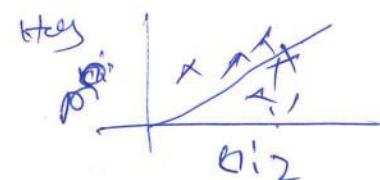
\rightarrow when $n = \mathbb{R}^n$ and $y = \mathbb{R}$, the simplest choice is to approximate y as a linear function of x

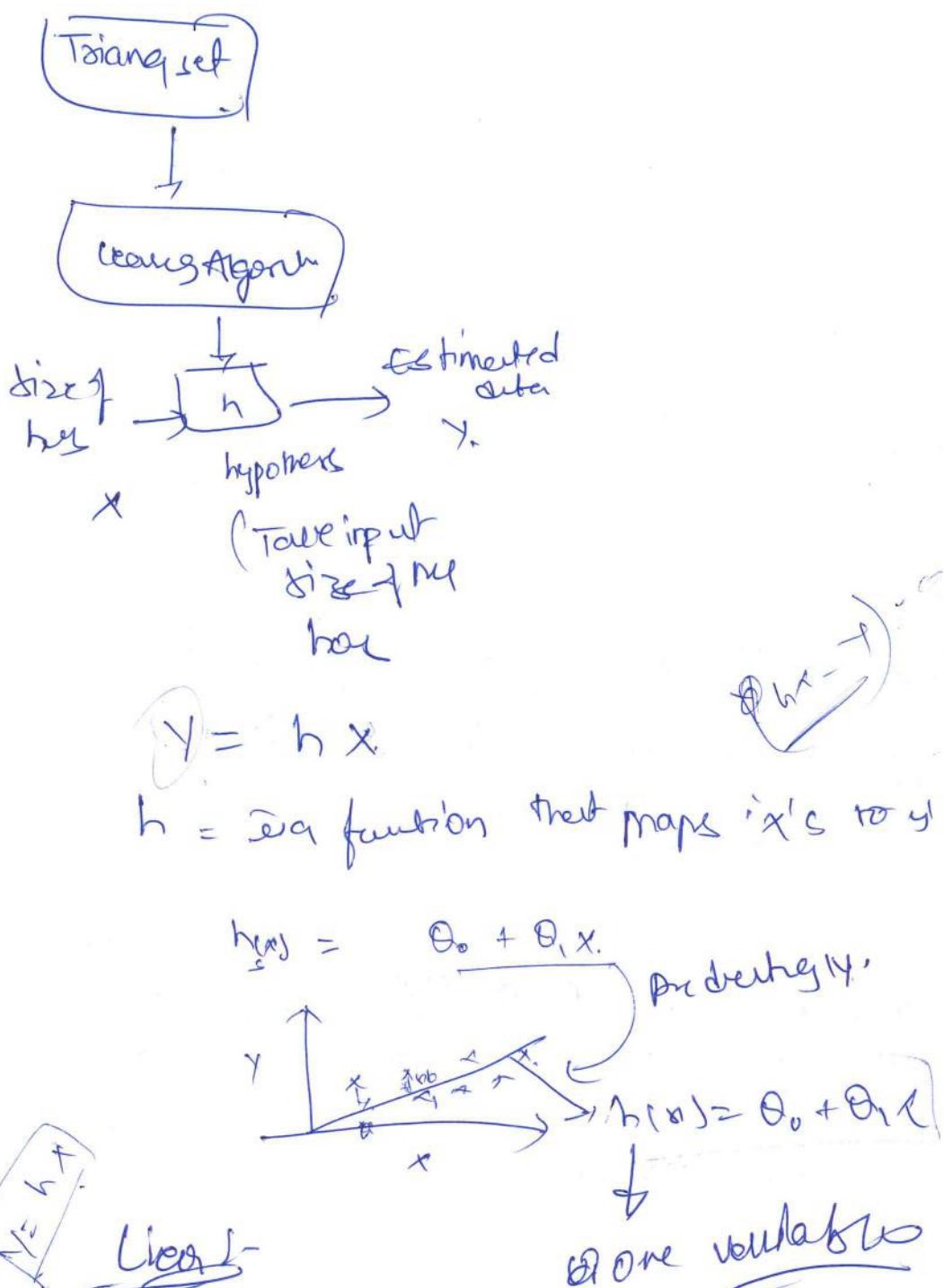
$$h_{\text{linear}} = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n.$$

x

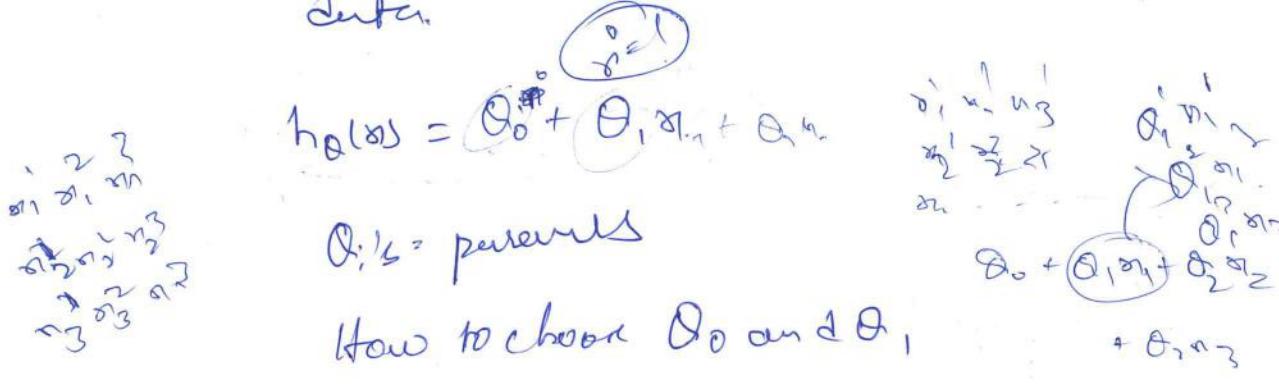
size x^n choice (h) (y)

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

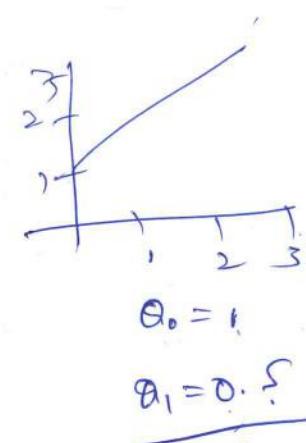
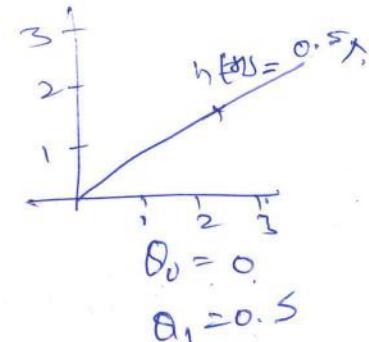
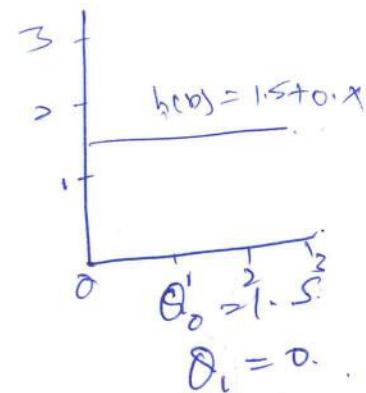




cost function: How to fit best possible straight line to our data.



How to choose θ_0 and θ_1



Idea: choose θ_0, θ_1 so that $h(x)$ is close to y for our training samples ($N \times bce$)

so minimize (θ_0, θ_1)

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

difference has to be minimum

$$\theta_0, \theta_1 \left[\sum_{i=1}^m [(h(\theta_0 + \theta_1 x_i) - y_i)^2] \right]$$

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

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

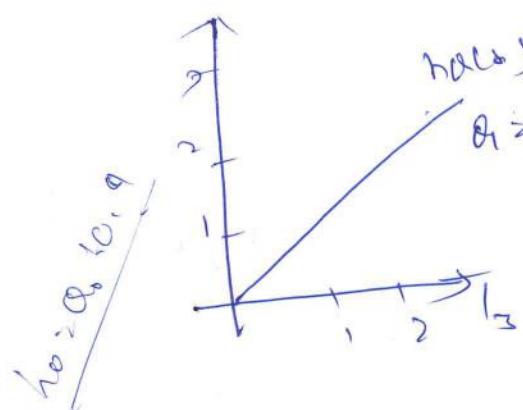
simplified

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



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

hole on



$$\theta_1 = 1$$

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

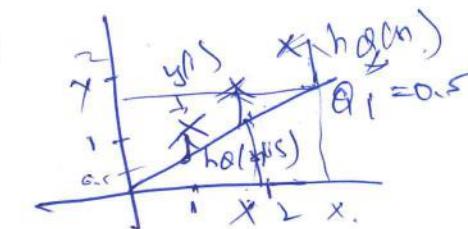
$$= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 =$$

constant for θ_1 and θ_2

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

$$J(\theta) = 0$$

If $\theta_1 = 0.5$



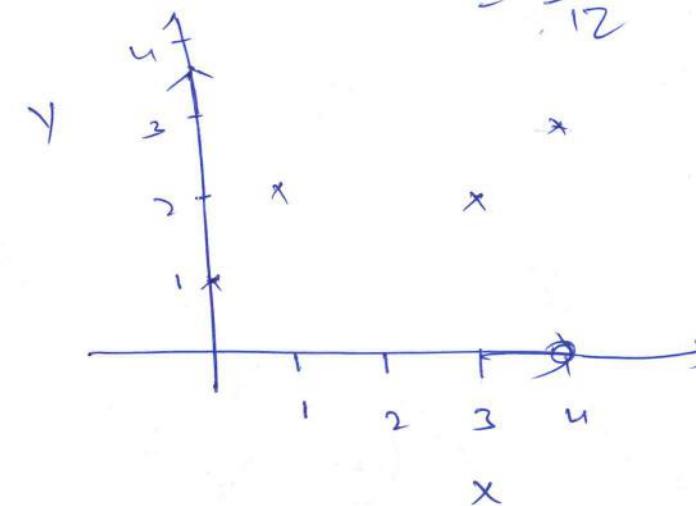
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$= m=3$$

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

$$= \frac{7}{6}$$

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



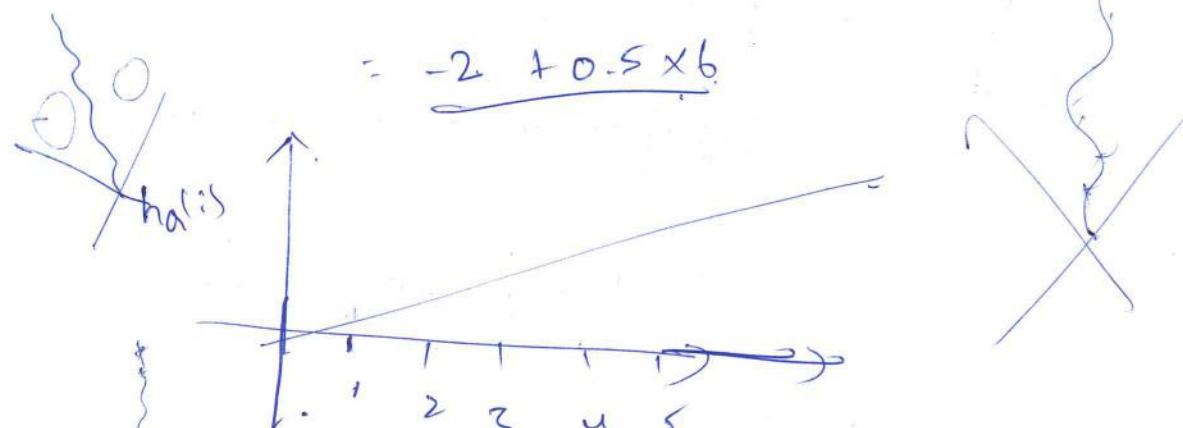
1700-D
cost = 0
 ≤ 0

$$-1780.0 + 11.5(520.9)$$

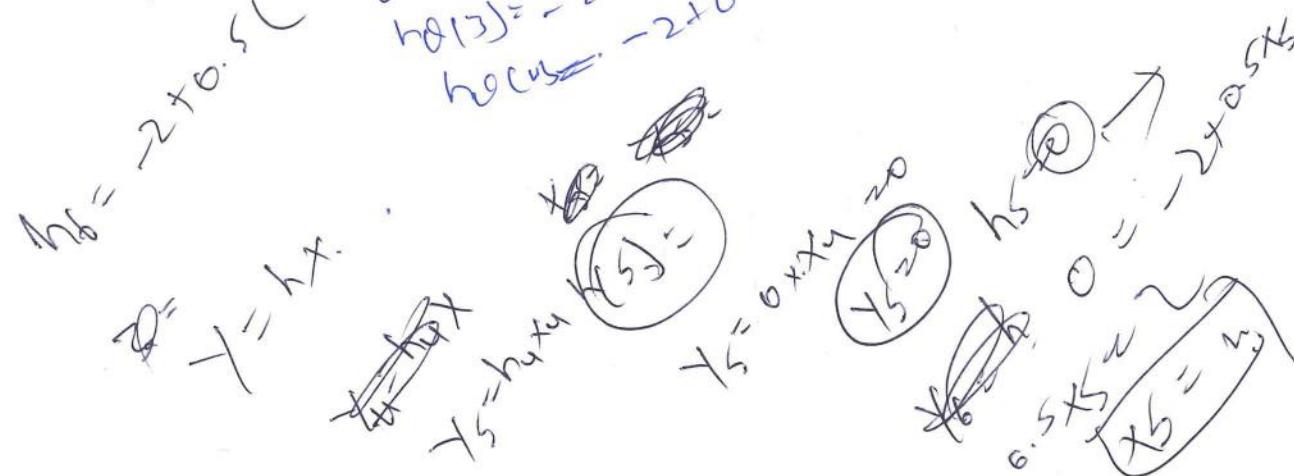
①

$$\text{h}_0(0) = \frac{1}{\theta_0} \{$$

$$\text{h}_0(6) = \theta_1 + \theta_2 \times 6 \quad \text{h}_0 = \theta_1 + \theta_2 x_0 \\ = -2 + 0.5 \times 6$$



$$\begin{aligned} \text{h}_0(0) &= -2 + 0.5(0) \Rightarrow -0.5 \\ \text{h}_0(2) &= -2 + 0.5(2) \Rightarrow -1.5 \\ \text{h}_0(4) &= -2 + 0.5(4) = -2 \\ \text{h}_0(6) &= -2 + 0.5(6) = 0 \end{aligned}$$



$$\text{h}_0(1) = -2 + 0.5(1) \Rightarrow$$

$$\text{h}_0(2) = -2 + 0.5(2) \Rightarrow$$

$$\text{h}_0(3) = -2 + 0.5(3) \Rightarrow$$

$$\text{h}_0(4) = -2 + 0.5(4) \Rightarrow$$

$$\text{h}_0(5) = -2 + 0.5(5) \Rightarrow$$

$$\text{h}_0(6) = -2 + 0.5(6) \Rightarrow$$

$$\text{h}_0(0)$$

$$\{ \text{h}_0(0), \text{h}_0(1), \text{h}_0(2), \text{h}_0(3), \text{h}_0(4), \text{h}_0(5), \text{h}_0(6) \}$$

② 4

$$\begin{aligned} 0 &= 1 \\ 3 &= 1 \\ n &= 1 \end{aligned}$$

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

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

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

$$\begin{aligned} \text{h}_0(2) &= 0 + 1 - 1(2) \\ &= 0 + 0 \end{aligned}$$

$$2 + 0 - 1(3) =$$

matrices :-

Wk 10 - 2

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

$$\begin{bmatrix} u \\ v \\ w \end{bmatrix} = \begin{bmatrix} 1 & 3 & 2 \\ 2 & 4 & 3 \\ 0 & 5 & 2 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 4 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 2 & 3 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 7 & 9 \\ 9 & 12 \\ 4 & 8 \end{bmatrix}$$

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

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

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

$$A = \begin{bmatrix} 1 & 4 \\ 2 & 2 \end{bmatrix}, C = \begin{bmatrix} 2 & 2 & n & 1 \end{bmatrix}$$

$$(A)(B)$$

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

$$v \cdot w = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 2+6+1 \\ 1 \end{bmatrix}$$

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

* Linear representation - multivariable.

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

In	In	In	Out
x_1	x_2	x_3	y
x_1^1	x_2^1	x_3^1	y^1
x_1^2	x_2^2	x_3^2	y^2
x_1^3	x_2^3	x_3^3	y^3

\Rightarrow multivariable.

n = number of features.

$\xrightarrow{\text{in}}$ input

$\xrightarrow{\text{out}}$ output

$x_{ij} = x_i^j$

$x_0 = 1$

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

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

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

$$h_\theta(x^3) = \theta_0 + \theta_1 (x_1^3) + \theta_2 (x_2^3) + \dots$$

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

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

$$= \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \theta_3 \end{bmatrix} \begin{bmatrix} x_0 & x_1 & x_2 & x_3 \\ x_1 & x_2 & x_3 & x_0 \\ x_2 & x_3 & x_0 & x_1 \\ x_3 & x_0 & x_1 & x_2 \end{bmatrix}$$

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

~~Octave~~

$$a = [1 \ 2 \ 3 \ 4 \ 5 \ 6]$$

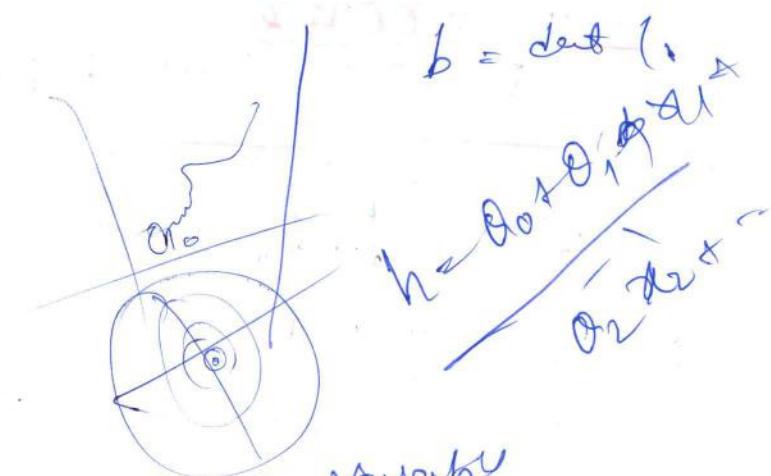
$$a = \begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix}$$

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

$$b = \begin{matrix} 1 \\ 4 \end{matrix}$$

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

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



gradient descent for multiple variables

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

↓
Update $\theta_0, \theta_1, \theta_2$

Update augman

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

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

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:

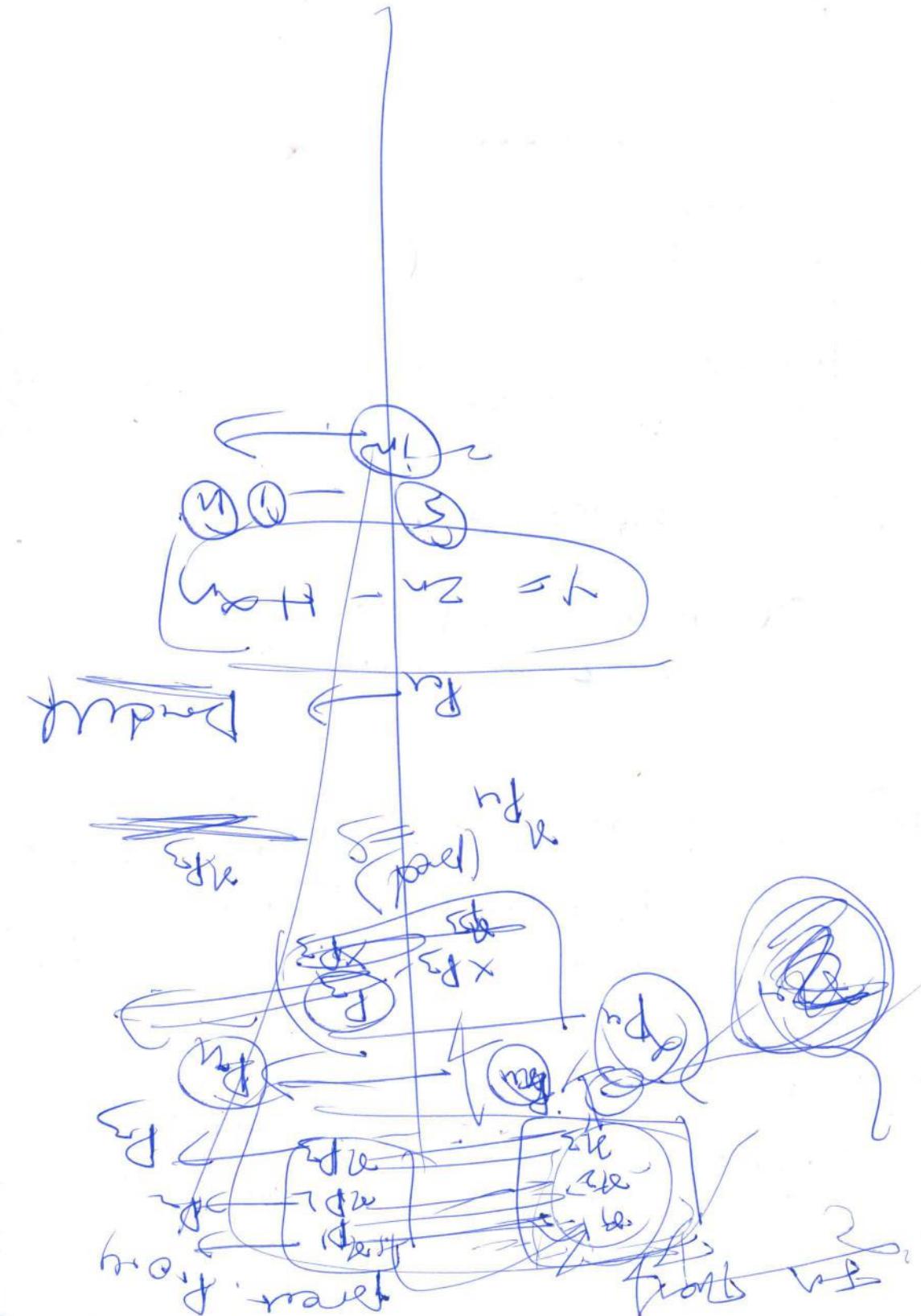
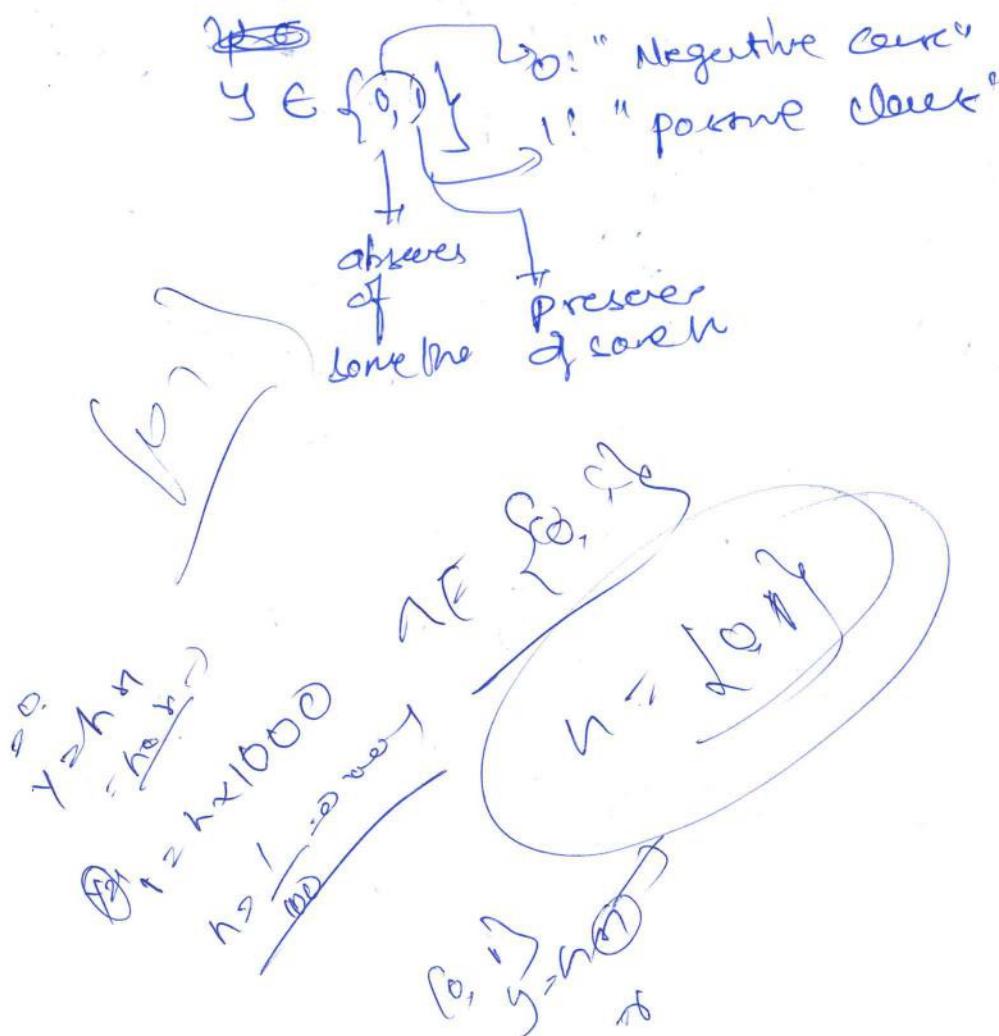
CLASSIFICATION AND Representation

Logistic Regression

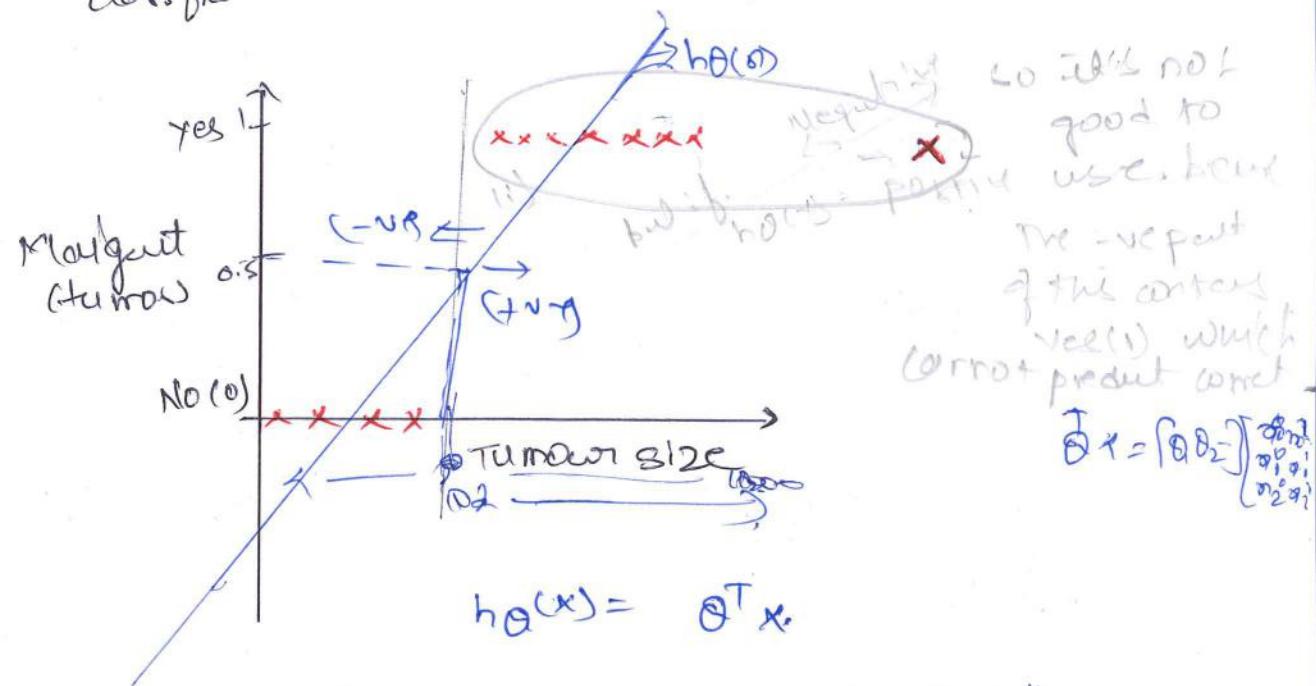
Classification :-

Error :- spam (Not spam)

online :- food item (Yes / No)



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



If $h_\theta(x) \geq 0.5$, predict "Y=1"
 $h_\theta(x) < 0.5$, predict "Y=0"

→ For Classification: $y \in \{0, 1\}$

$h_\theta(x)$ can be > 1 or < 0

so error
not good.

→ so, we go to

Logistic regression: $0 \leq h_\theta(x) \leq 1$
also rate of classifier

Logistic Regression: Hypothesis Regression.

want $0 \leq h_\theta(x) \leq 1$

so

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

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

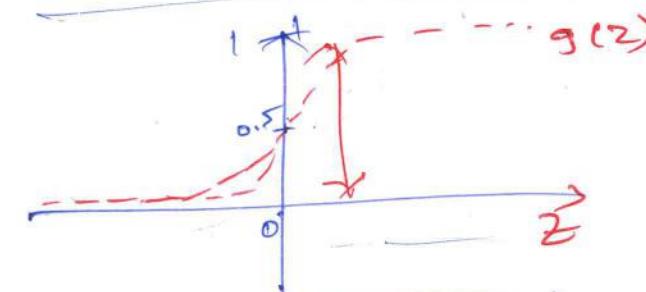
$$= g(z)$$

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

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

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} = \frac{1}{1+e^{-\theta_0 - \theta_1 x_1 - \theta_2 x_2}}$$

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

(8) 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 classification

$y \geq 0$ or 1.

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

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

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

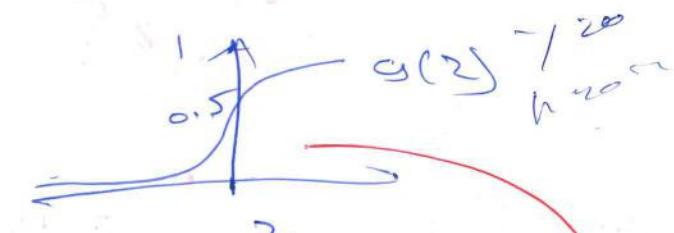
~~$P(y=0/x; \theta) = 1 - P(y=1/x; \theta) = 1 - 0.7 = 0.3$~~

Decision Boundary

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

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

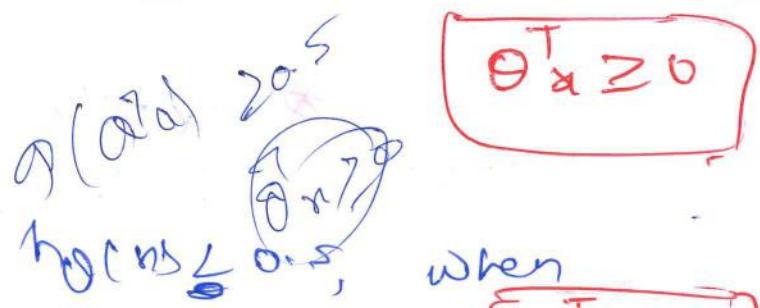
when it makes $y=1$ and $y=0$ easier to understand the hypothesis (particularly when there are many features).



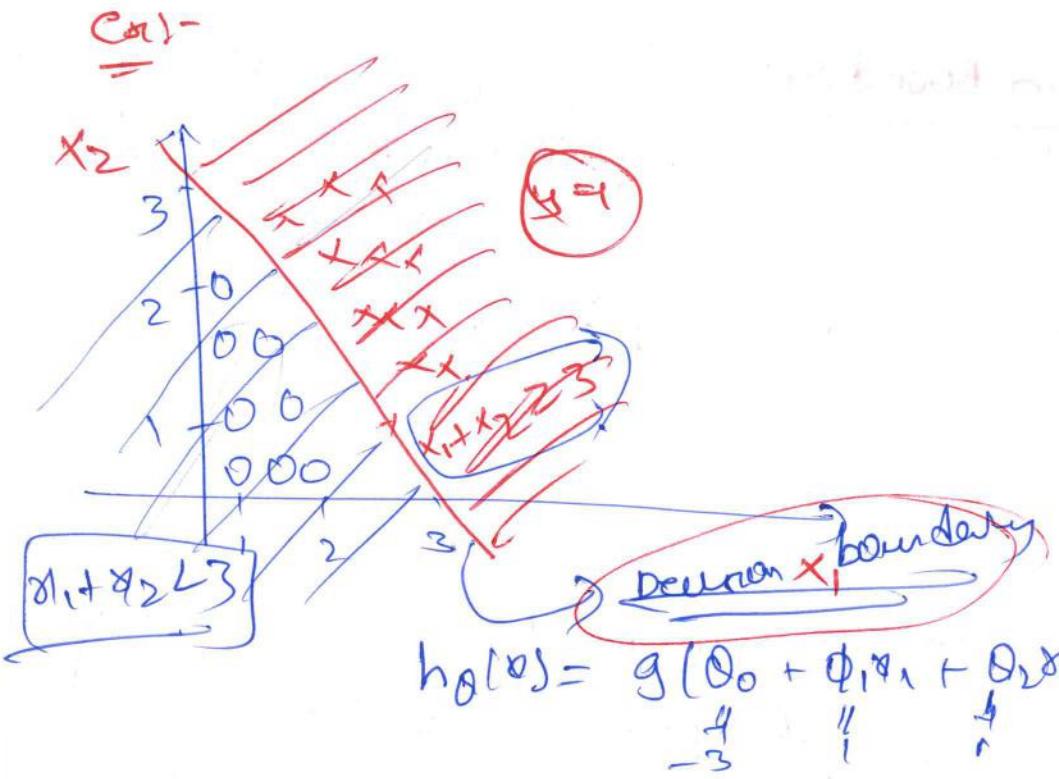
Suppose predict " $y=1$ " if $h_{\theta}(x) \geq 0.5$
predict " $y=0$ " if $h_{\theta}(x) < 0.5$

→ To understand better, we need when
 ~~$y \geq 0$ given $h_{\theta}(x) \geq 0.5$~~
when $g(z) \geq 0.5$
when $z \geq 0$

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



→ Similarly $h_{\theta}(x) \leq 0.5$, when $\theta^T x \leq 0$



Predict $\hat{y} = 1$, if $(-3 + \theta_1 + \theta_2 \geq 0)$

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

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

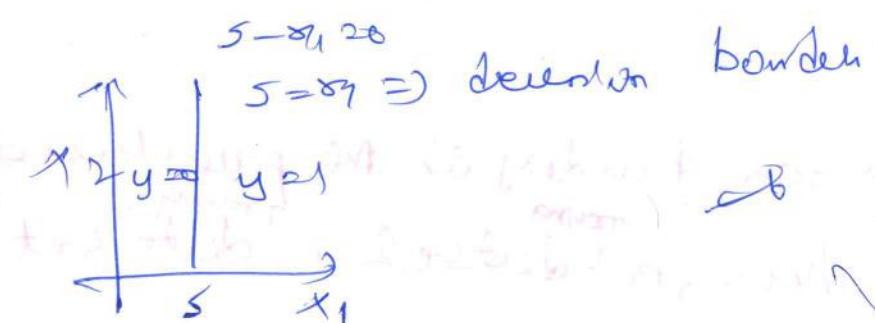
$$x_1 + x_2 < 3 \Rightarrow y = 0.$$

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

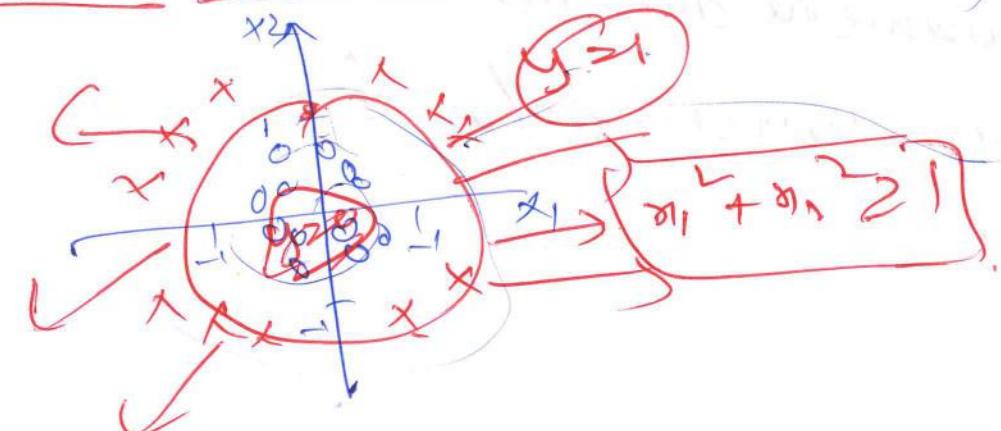
Cont-

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

$$\theta_0 = 5, \theta_1 = -1, \theta_2 = 1$$



Non linear decision boundaries



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

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

Predict y' if $-1 + x_1^2 + x_2^2 \geq 0$

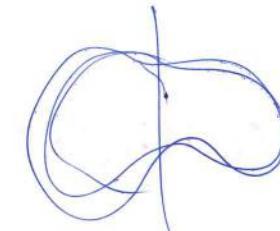
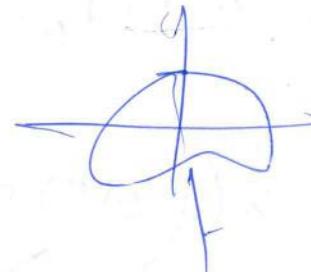
$$x_1^2 + x_2^2 \geq 1 \Rightarrow \begin{cases} x_1 \geq 0, x_2 \geq 1 \\ x_2 \geq 0, x_1 \geq 1 \end{cases}$$

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

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

* more complex decision boundaries

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



Logistic Regression: cost function! How to

fit parameter into to fit to training set

ex: Training set = $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}) \dots (x^{(m)}, y^{(m)})\}$

example no $\{x^{(1)}, y^{(1)}\}, n_0 = 1, y_0 = 0, y_1 = 1$

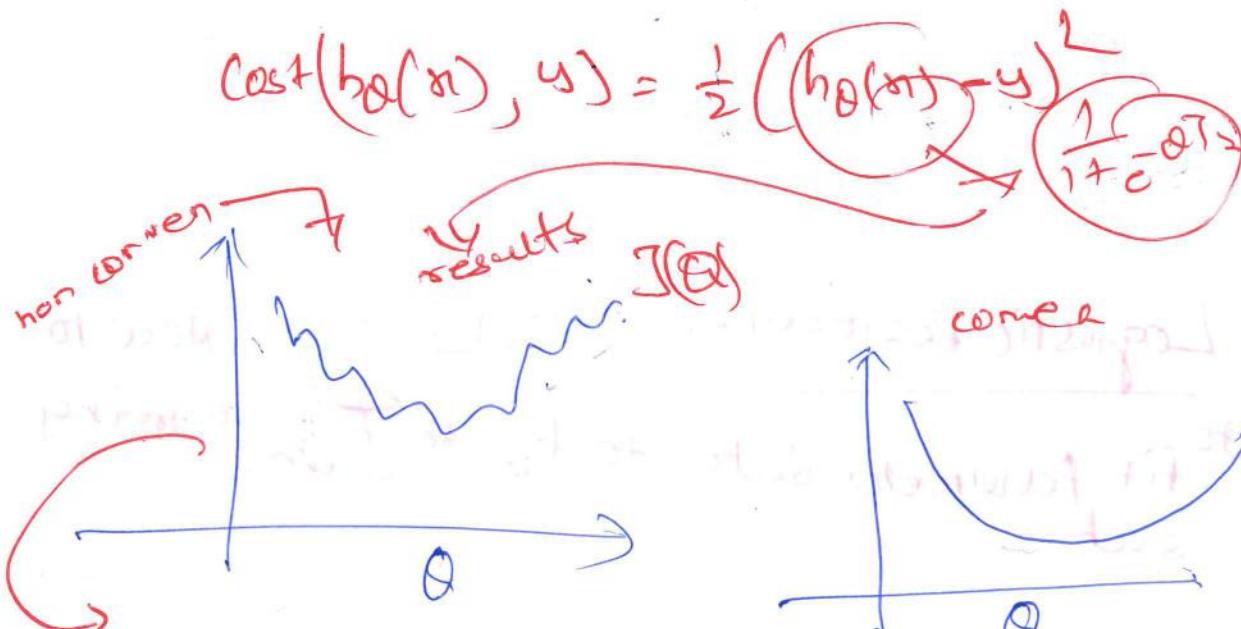
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose θ ?

cost function

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

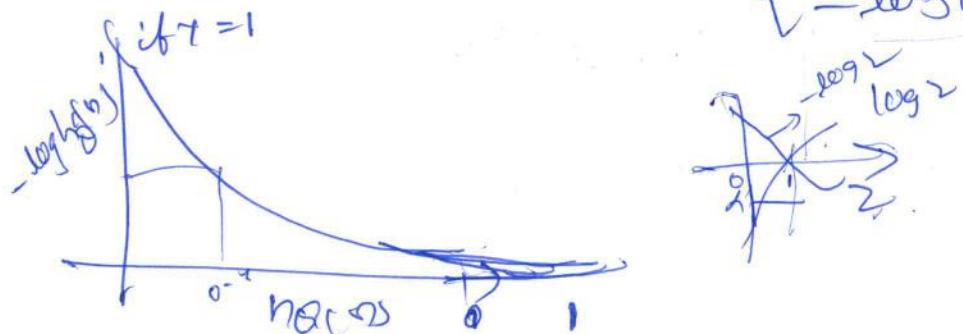
Logistic Regressn. $\Rightarrow \text{cost}(h_{\theta}(x^{(i)}), y^{(i)})$



\rightarrow * SO we cannot use the same cost function for logistic, in order to avoid many local minima.

Logistic regression cost function \rightarrow so use

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



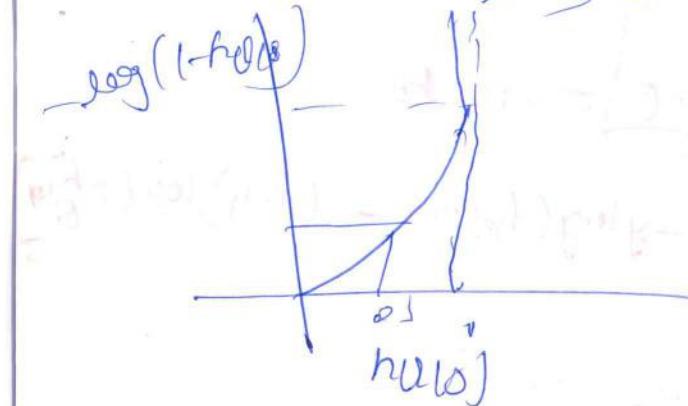
$\rightarrow \text{lost} = 0, \text{ if } y = 1, h_{\theta}(x) = 1$

But as $h_{\theta}(x) \rightarrow 0$
cost $\rightarrow \infty$

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

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

\rightarrow If $y = 0 \rightarrow$ large cost.



Sum

$$\begin{aligned} &-\log(h_{\theta}(x^{(i)})) \cdot y^{(i)} - \log(1 - h_{\theta}(x^{(i)})) \cdot (1 - y^{(i)}) \\ &- \log(1 - h_{\theta}(x^{(i)})) \cdot 2y^{(i)} - \log(h_{\theta}(x^{(i)})) \cdot (2y^{(i)} - 1) \\ &- \log(h_{\theta}(x^{(i)})) \cdot (1 - y^{(i)}) - \log(1 - h_{\theta}(x^{(i)})) \cdot (y^{(i)} - 1) \end{aligned}$$

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 criterion

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

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

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

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

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

To fit parameters θ :

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

⇒ To make a prediction given new x :

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

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 scope

Example:

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}, J(\theta) = (\theta_1 - s)^2 + (\theta_2 - t)^2$$

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

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

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

octave example

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

for loop

function [jVal, gradient] = costFunction(theta)

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

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

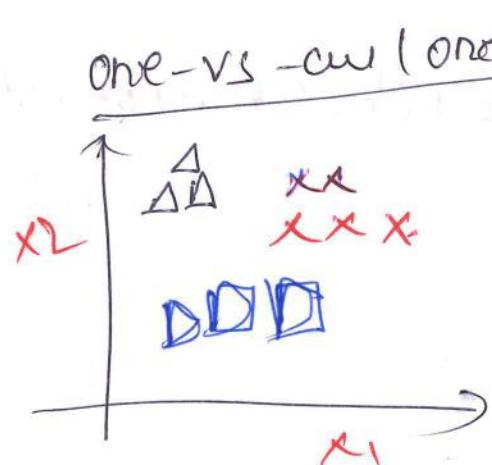
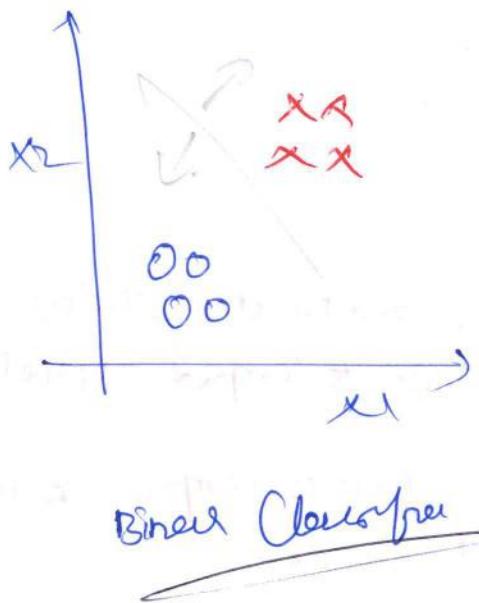
end.

Logistic Regression: MULTI CLASSIFICATION | One-vs-all.

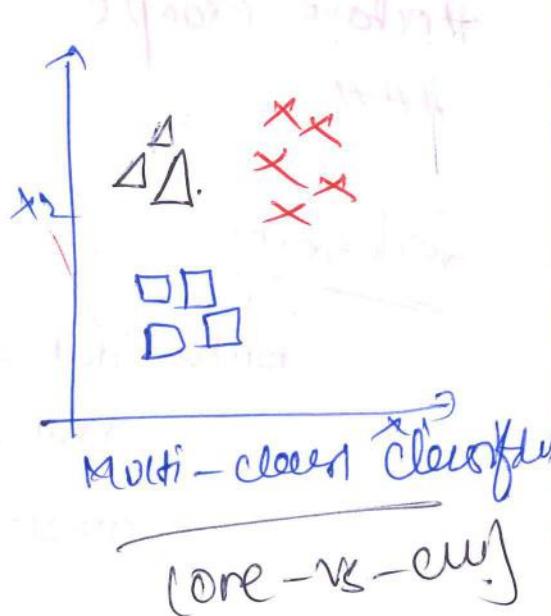
Example: work, friends, family, Hobby

Medical diagnosis: ~~Not ill~~, cold, flu

Weather: sunny, cloudy, rain, snow



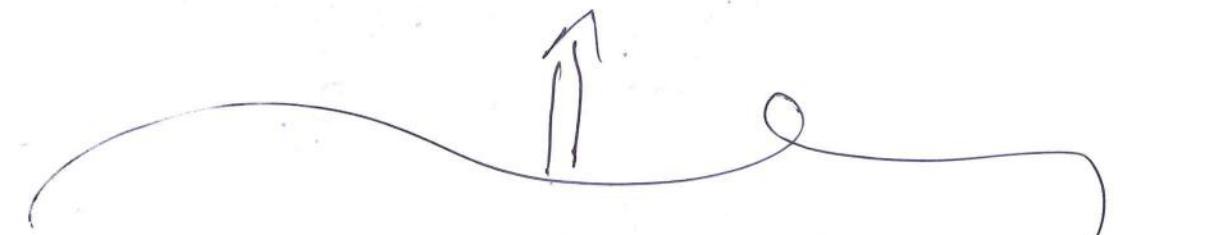
Class 1: A
Class 2: B
Class 3: C



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

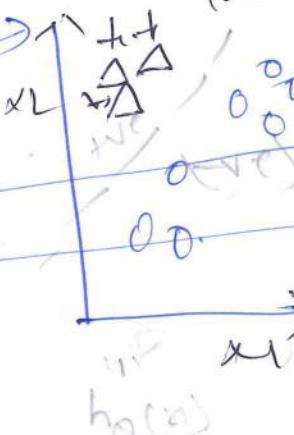
one-vs-one classifier has for every class \hat{y}^{cls}
to predict the probability that $y = \hat{y}^{\text{cls}}$

On a new input x , to make a prediction, we
choose that maximises $\max_i h_{\theta}^{\text{cls}}(x)$.



class 0

one-vs-all
(one-vs-great)



class 1

one-vs-all

(one-vs-great)



class 2



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

$$h = \theta^T x$$

$$\begin{bmatrix} h_0 \\ h_1 \\ h_2 \end{bmatrix} = \theta \begin{bmatrix} \theta_0, \theta_1, \theta_2 \end{bmatrix} \begin{bmatrix} x_0 & x_1 & x_2 \\ n_0 & n_1 & n_2 \\ n_0^2 & n_1^2 & n_2^2 \\ \dots & \dots & \dots \\ n_0^3 & n_1^3 & n_2^3 \end{bmatrix}$$

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

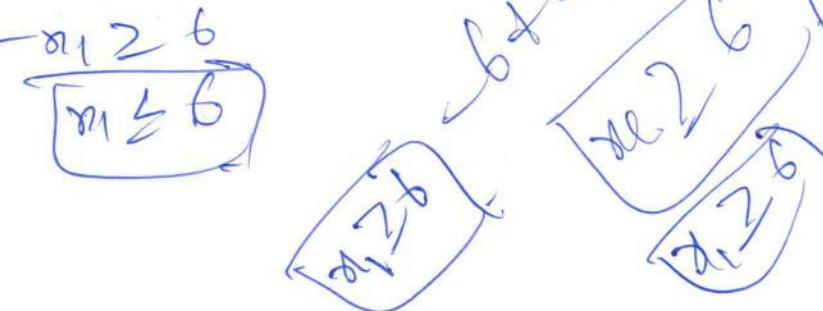
$$\cos(h - y) =$$

$$b + x_1 \geq 0$$

$$x_1 \geq 0$$

$$b + x_1 \leq 1$$

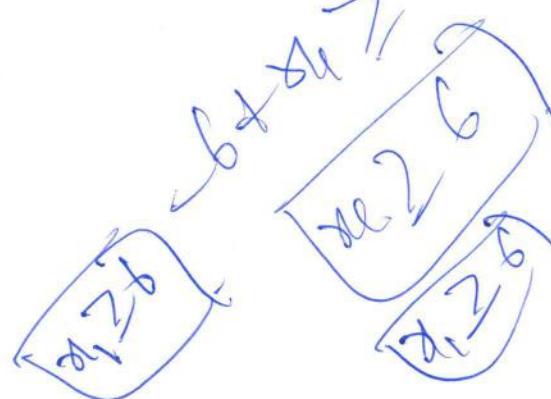
$$x_1 \geq 0$$



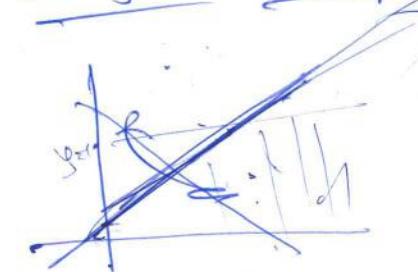
* $h = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$

$$b + x_1 \geq 0$$

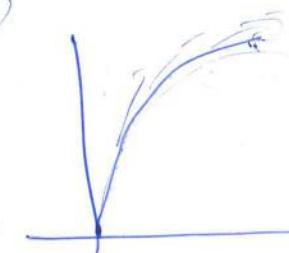
$$-x_1 \geq b$$



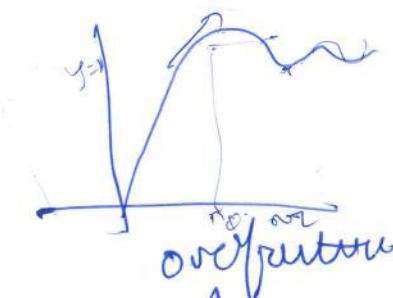
* Dangers of overfitting



underfit



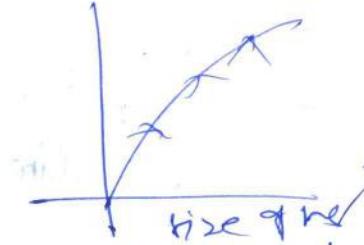
OK



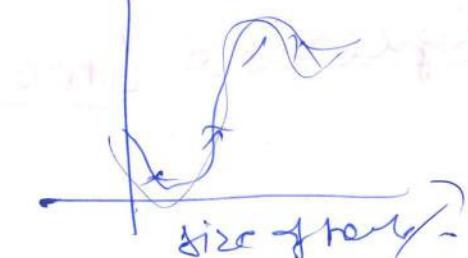
overfit

Can be avoided by choosing E
the no of parameters.

cost function



$$\theta_0 + \theta_1 x_1 + \theta_2 x_2$$



$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

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

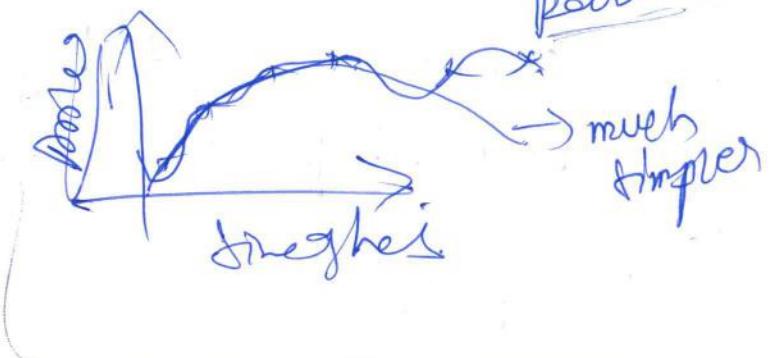
$$\theta_3 \approx 0, \theta_4 \approx 0$$

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

less overfit

degreenized :-

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



gegenseitige
Rechte

reeta

Much
Hunger

of Required lines

→ Grand desert
über seges.

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

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

② Crochet desert
litter bag

$$\hat{Q}_0 := Q_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(\hat{g}_0(\hat{x}_i^{(t)}) - y_i^{(t)} \right)$$

$$\rightarrow \theta_j = \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n h(\tilde{x}_i)$$

$$= y^{(1)} \quad \text{at } y^{(1)}$$

$$j = \cancel{0}, 1, 2, 3 \dots n$$

beginning

$$\theta_j = \theta_j - \alpha \left(\left[\frac{1}{m} \sum_i (\text{label}(i) - y^{(i)}) \right] \eta_j^{(i)} + \frac{\lambda}{m} \theta_j \right)$$

$$\left[\frac{1}{m} \sum_i (h(x_i) - y_i^{(j)}) w_j^{(i)} + \frac{2}{m} \theta_j \right]$$

$$j=1$$

Request for JCO

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

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

$\frac{1}{2} \pi r^2$

$$\text{Qj} \propto 0.9$$

enrich
by oil

No need to queue

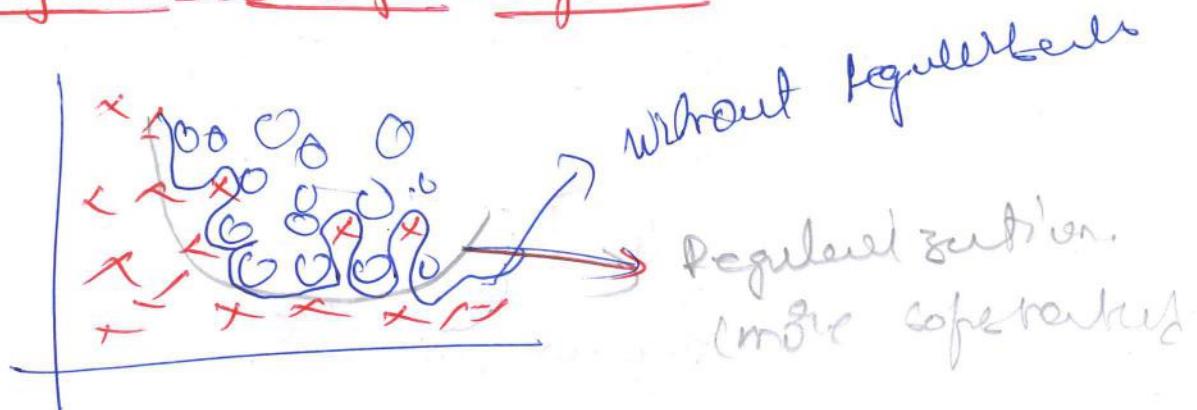
$$X = \begin{bmatrix} v \\ w \\ \vdots \\ z^{(m)} \end{bmatrix}^T$$

$$y = \begin{bmatrix} y^{(0)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

$$\min j(\phi)$$

$$\alpha = (x^T x)^{-1} \lambda \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} x^T$$

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

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

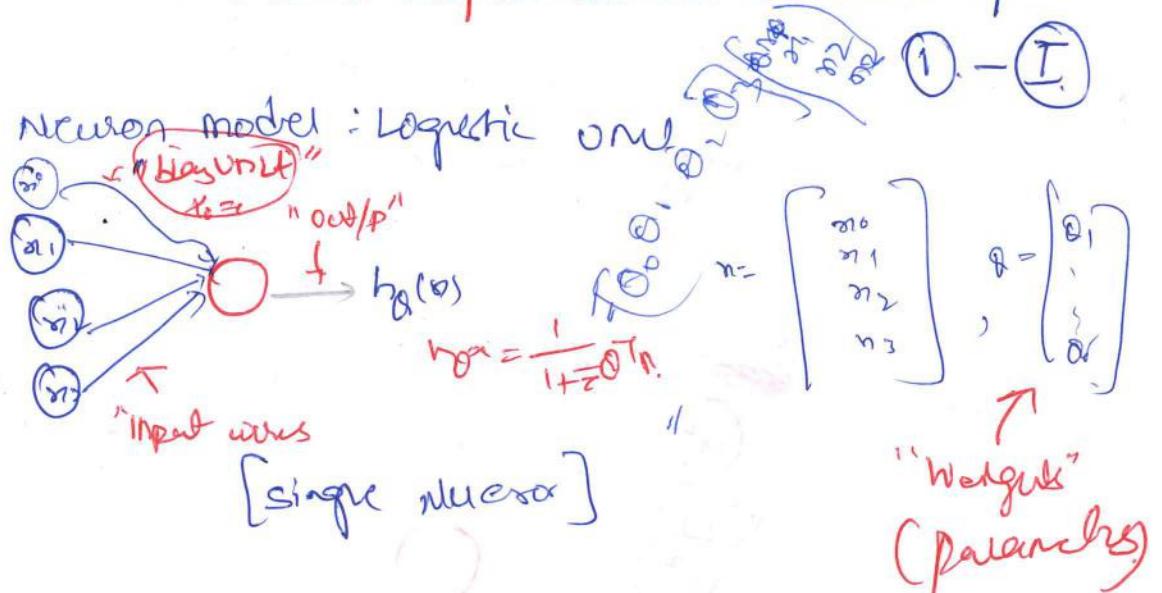
$$\theta_j := \theta_j - \alpha \left\{ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right\}$$

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

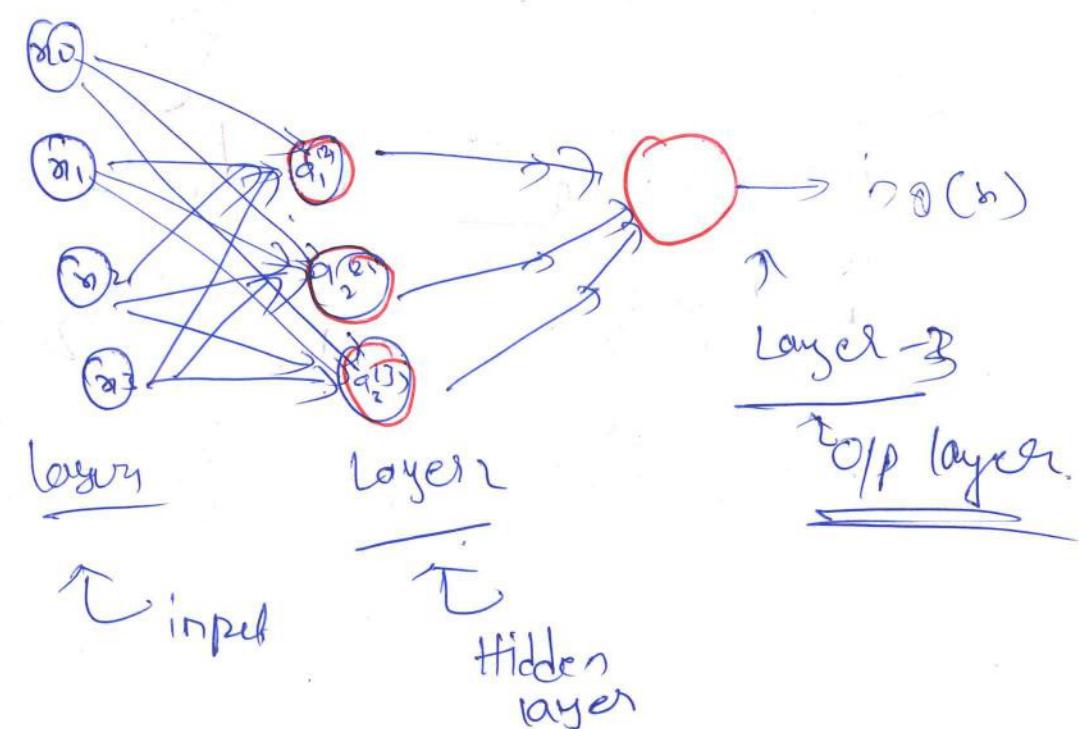
WEEK 4

NEURAL NETWORK: REPRESENTATION

— Model Representation: Model representation

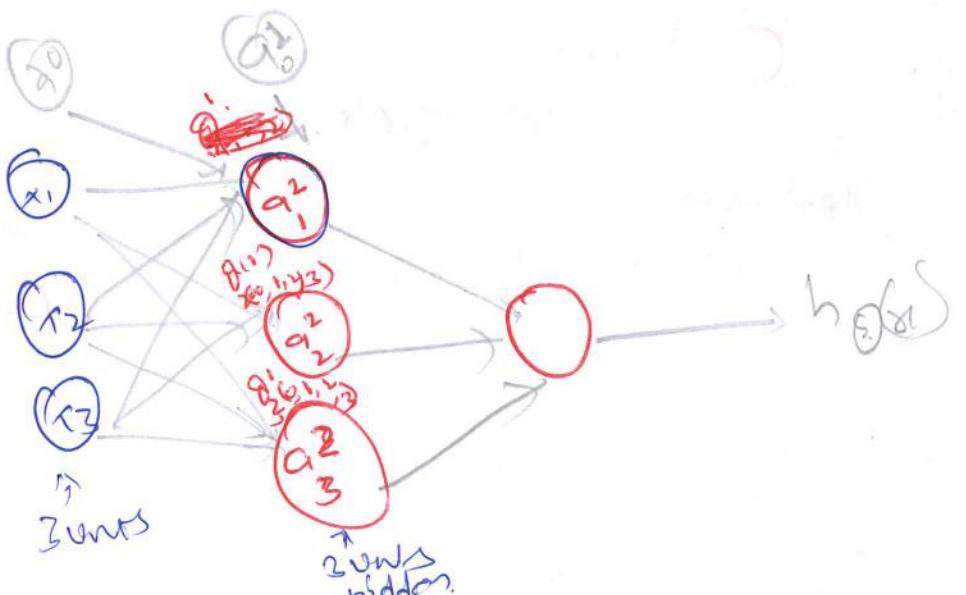


Neural Network (multiple)



$a_i^{(j)}$ = "activation" of unit i in layer j

$\Theta^{(j)}$ = matrix of weight connecting function
maps from layer j to layer $j+1$



$$a_1^{(2)} = g(\theta_{00}^{(2)} x_0 + \theta_{10}^{(2)} x_1 + \theta_{20}^{(2)} x_2 + \theta_{30}^{(2)} x_3)$$

$$a_2^{(2)} = g(\theta_{01}^{(2)} x_0 + \theta_{11}^{(2)} x_1 + \theta_{21}^{(2)} x_2 + \theta_{31}^{(2)} x_3)$$

$$a_3^{(2)} = g(\theta_{02}^{(2)} x_0 + \theta_{12}^{(2)} x_1 + \theta_{22}^{(2)} x_2 + \theta_{32}^{(2)} x_3)$$

If in layer $j = s_j$ units

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

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

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

↓
Gute Rechner.

$$h_{\Theta}(x) = a^{(2)} = g(\theta_{00}^{(2)} x_0 + \theta_{10}^{(2)} x_1 + \theta_{20}^{(2)} x_2 + \theta_{30}^{(2)} x_3)$$

$$+ g(\theta_{01}^{(2)} x_0 + \theta_{11}^{(2)} x_1 + \theta_{21}^{(2)} x_2 + \theta_{31}^{(2)} x_3)$$

$$+ g(\theta_{02}^{(2)} x_0 + \theta_{12}^{(2)} x_1 + \theta_{22}^{(2)} x_2 + \theta_{32}^{(2)} x_3)$$

Model representation + U

Consider here,

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

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

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

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

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

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

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

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

Add $a_0 = 1$

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

$$h_\theta(z) = a^3 = g(z^3)$$

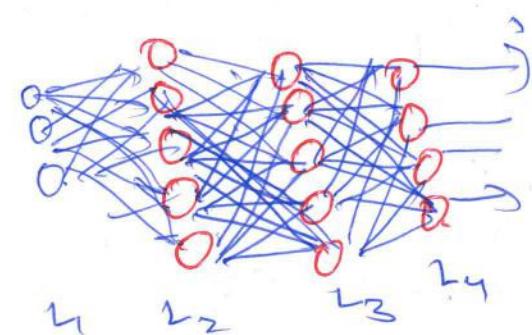
WEEK: 5

 =

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

= n
 $s_L = \text{no. of units}$
(not counting bias units)

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

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

multiple output

$$h_\theta(x) \in \mathbb{R}^{1C}$$

$1C = \text{no. of o/p.}$

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

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 $l^{(i-1)}$

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

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

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

$$\begin{matrix} 0 \\ 0 \\ 0 \end{matrix} \rightarrow \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} \rightarrow \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} \quad y^{(0)} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

No regularization of bias terms only terms $\theta_1, \theta_2, \dots, \theta_n$

BACK PROPAGATION: ~~MIN~~ MIN cost function

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

↳ need to compute $-J(\theta)$

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

$i = \text{No. of features of each } x^{(i)}$

$j = \text{No. of e}$

number

$i = \text{No. of each } x_j = \theta_j^1, \theta_j^2, \dots$
 $= (car, bus, car, \dots)$

$j = \text{No. of patterns} = (\text{car, bus, toy})$

Gradient computation
 Given on training example ($n \times Y$)
 not $\theta_1, \theta_2, \dots, \theta_n$

$$\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \rightarrow \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \rightarrow \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \quad \theta^{(3)} \rightarrow \text{ans}$$

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

for one
(x, y)

: Back propagation algorithm

→ Intuitively consider an error, $\delta_j^{(l)}$ = error

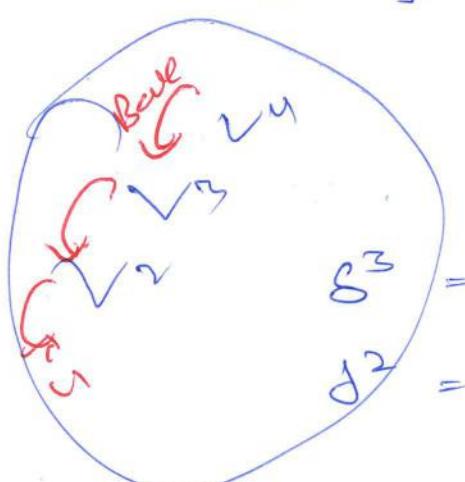
of node j in layer l .

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

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

$$h_{\Theta}(x, a_j)$$

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



$$\delta^3$$

$$\delta^3 = (\Theta^{(2)})^T \delta^{(2)} \cdot g'(z^{(2)})$$

$$\delta^2 = (\Theta^{(1)})^T \delta^{(1)} \cdot g'(z^{(1)})$$

$\delta^{(1)}$ of first node

$$\frac{\partial J(\Theta)}{\partial \theta_{ij}} = g_j^{(l)} s_i^{(l+1)}$$

(ignore s_i if $i=0$)

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)

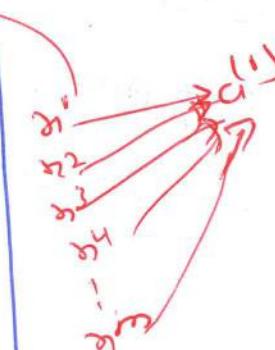
(used to compute

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

$$\text{set } \delta^{(L)} = \delta^{(L)}$$

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

$$\frac{\partial}{\partial \theta_{ij}} J(\Theta)$$



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

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

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

Writing
Aggs in fresh : Backpropagation algm

→ 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 $a^1 = s^{(1)}$ ~~for $\forall i$~~

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

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

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

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

→ done after this

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

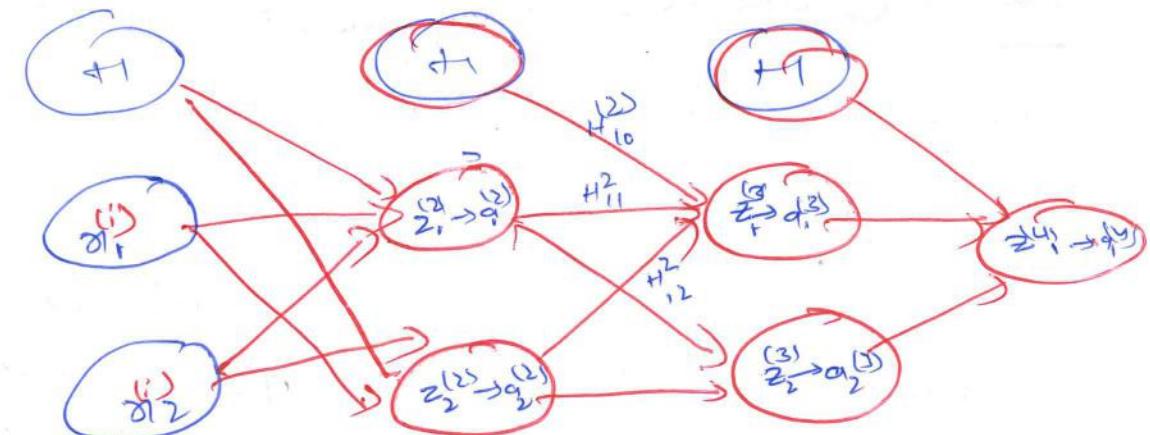
$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \left. \begin{array}{c} \frac{\partial J(\Theta)}{\partial \Theta_{ij}^{(l)}} \\ \Delta_{ij}^{(l)} \end{array} \right\} \left. \begin{array}{c} D_{ij}^{(l)} \\ \Delta_{ij}^{(l)} \end{array} \right\}$$

Out

Out of Context:

Forward Propagation [input units]

Understands Back Propagation



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

$$z_1^{(3)} = \Theta_0^{(3)} x_1 + \Theta_1^{(3)} a_1^{(2)} + \Theta_2^{(3)} a_2^{(2)}$$

→ In general $s_{ij}^{(l)} = \text{"error of cost. } a_{ij}^{(l)} \text{ - (unit } i \text{ in layer } l \text{) }$

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

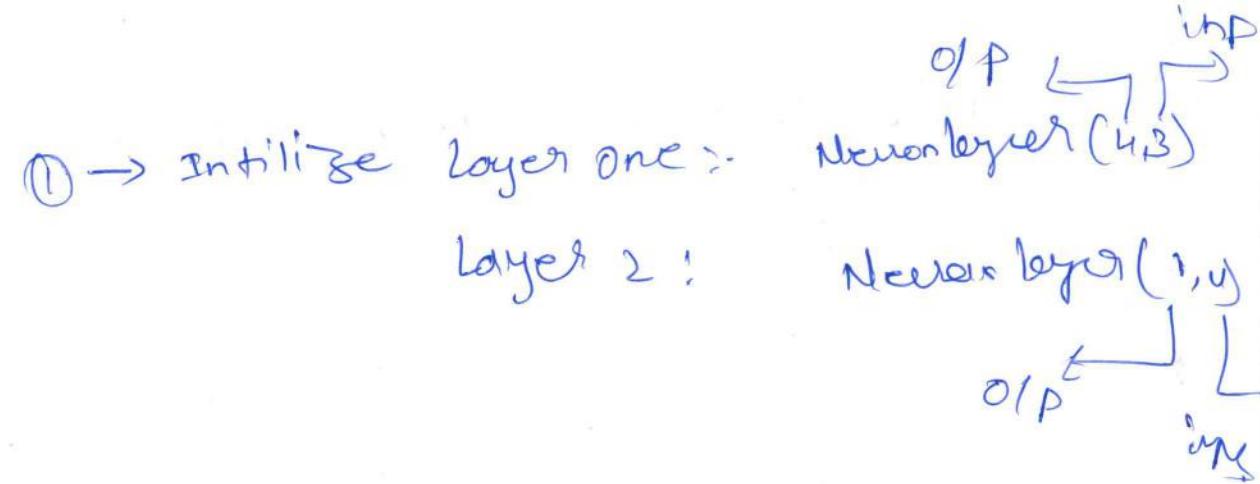
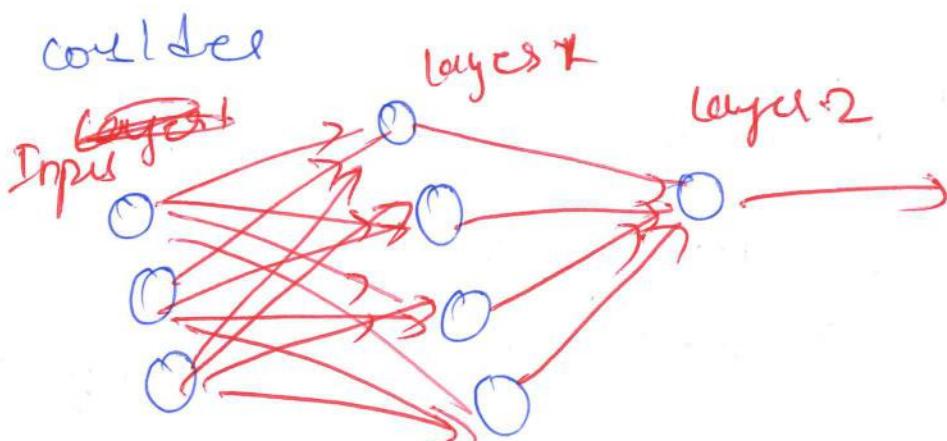
→ $\text{cost}(i) = y^{(i)} \log h_{\Theta}(a^{(i)}) + (1-y^{(i)}) \log h_{\Theta}(a^{(i)})$ where

Back propagation In Python

Implementation note: unrolling parameters.

Out of Course: Implementing multilayer neural networks. In general single layer causes error when there is large input data in which it cannot generate to me. O/P having size

so.
steps for implementing



① → Initialize layer one: Neuron layer (4,3)
Layer 2: Neuron layer (1,1)

- ② → Initialize training set 'matrix'
- ③ → Initialize training set output matrix
- ④ → Initialize random weights
- ⑤ → ~~Find~~ find O/P → find error b/w actual L & O/P and obtain error (crossing)
- ⑥ → ~~B~~ feed this ans
→ find O/P from L₂ by feed day L₁ input to L₂

- Ⓐ → find real L_2 and $\partial D L_2$
cross
- Ⓑ → Use error from $-L_2'$ time
 ΔL_2 , then ~~for~~ ΔL_2
Can be found using sign & 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

→ ~~fixe ~~high~~ ~~high~~ areas~~
By (high var)

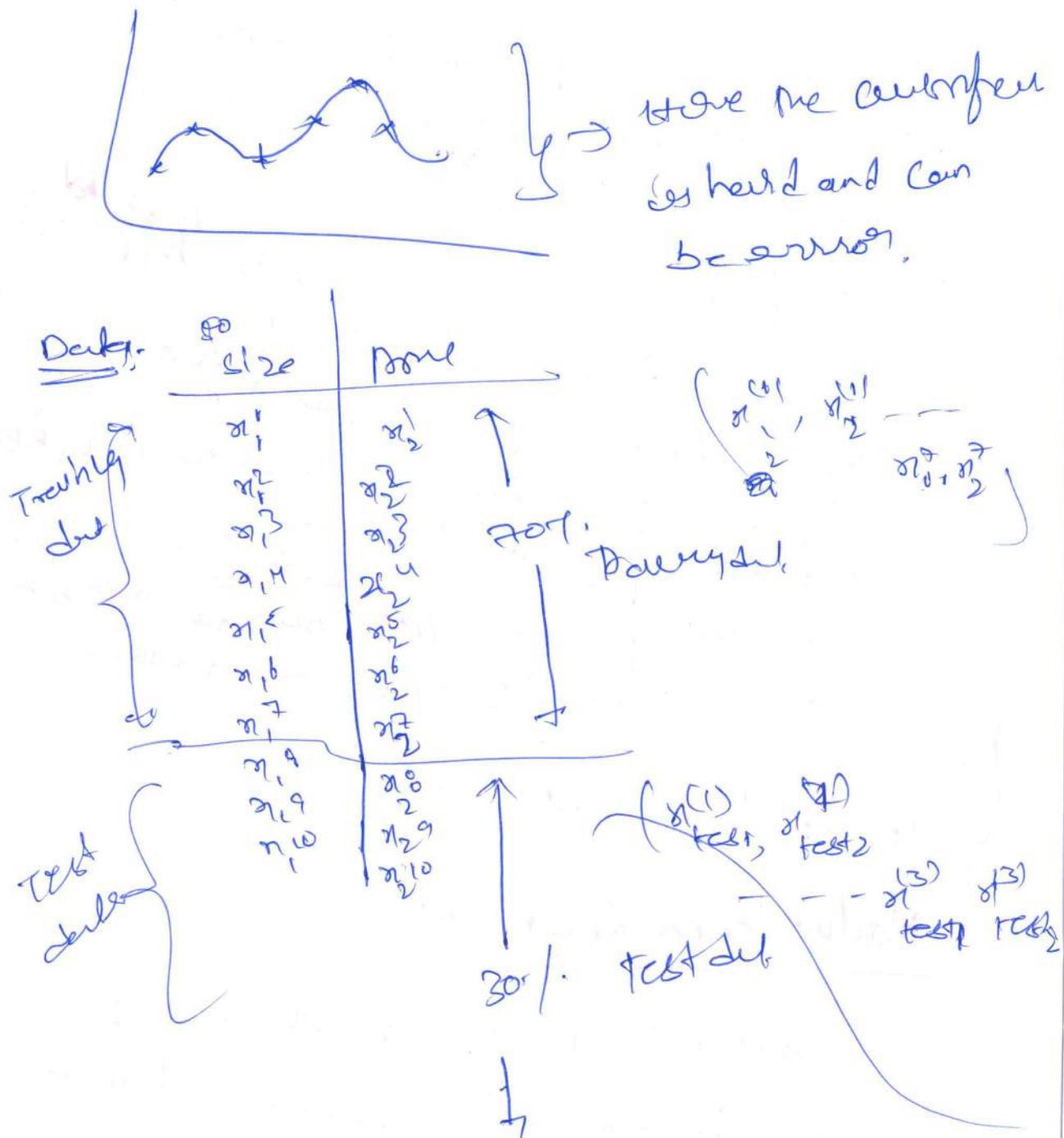
- fix more training set. ✗ (bad idea)
- Try smaller sets of features → ~~high by high variance~~
- try getting additional features → fixed by ~~high~~
- Try adding polynomial features ($x^1, x^2, x^n, \dots, e^{10}$) → fixes ~~high~~ ~~high bias~~
- Try deleting λ → fixes ~~high~~ ~~high bias~~
- Try increasing λ → fixes ~~high~~ ~~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



Have the overfitted
is hard and can
be error.

↑ 60% training
↓ 10% cross validation
↓ 20% test

\Rightarrow

Training error:

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

Cross validation error:

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

Test error

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

\Rightarrow Model selection

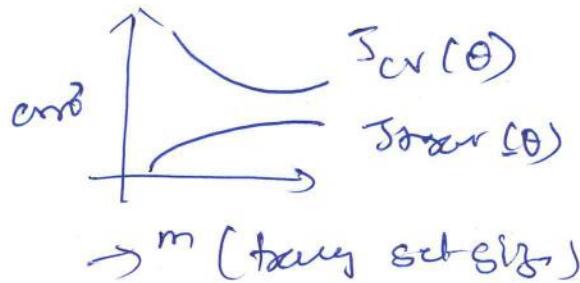
1. $h_{\theta}(x) = \theta_0 + \theta_1 x_1$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$

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

→ APPLY $J_{test}(\theta)$

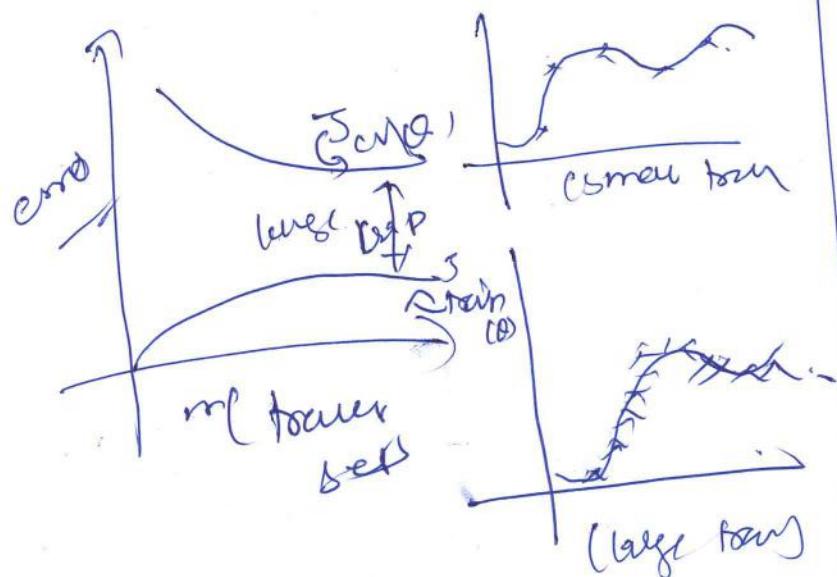
Learned levels :-

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



For high bias

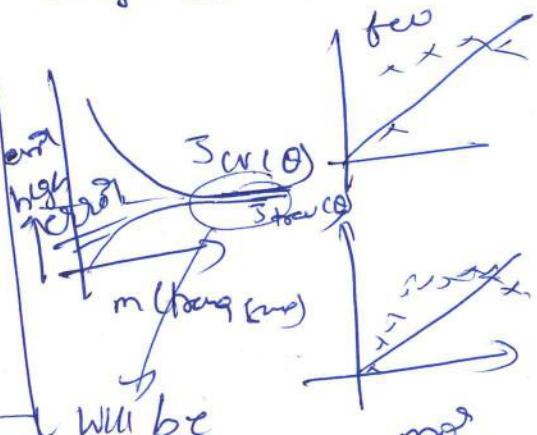
$$h_{\theta \text{ bias}} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots - \theta_{100}^{100}$$



For high bias

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

involves lot more training



will be close

high bias, high error will be
 $J_{cv}(\theta)$, J_{train} almost same.

Small Network
Network

computationally
cheaper

Large network
Network

computationally
more expensive

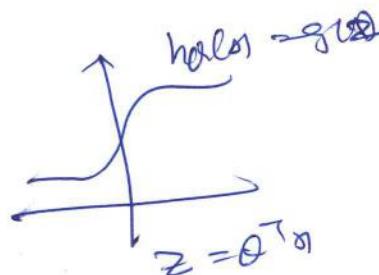
use regularization (λ)
to address overfitting,

SUPPORT VECTOR MACHINES

I) OPTIMIZATION objective :-

Alternate view of logistic regression

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

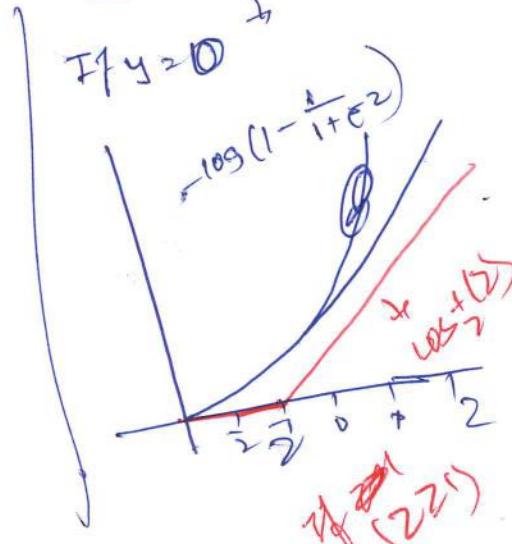
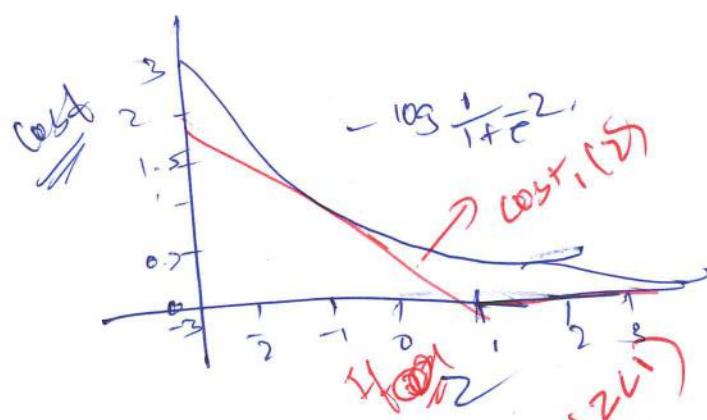


If $y=1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x \gg 0$
 $y=0$.. $h \approx 0$. $\theta^T x \ll 0$

* 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:

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

so

SVM Hypothesis:

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

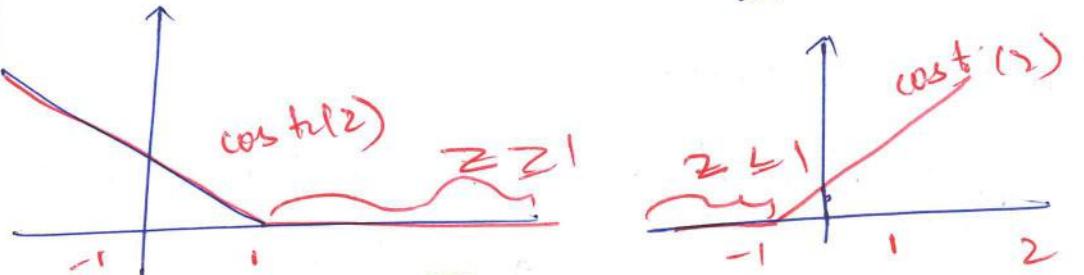
Hypothesis:

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

SN

SVM: Large Margin Classification

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



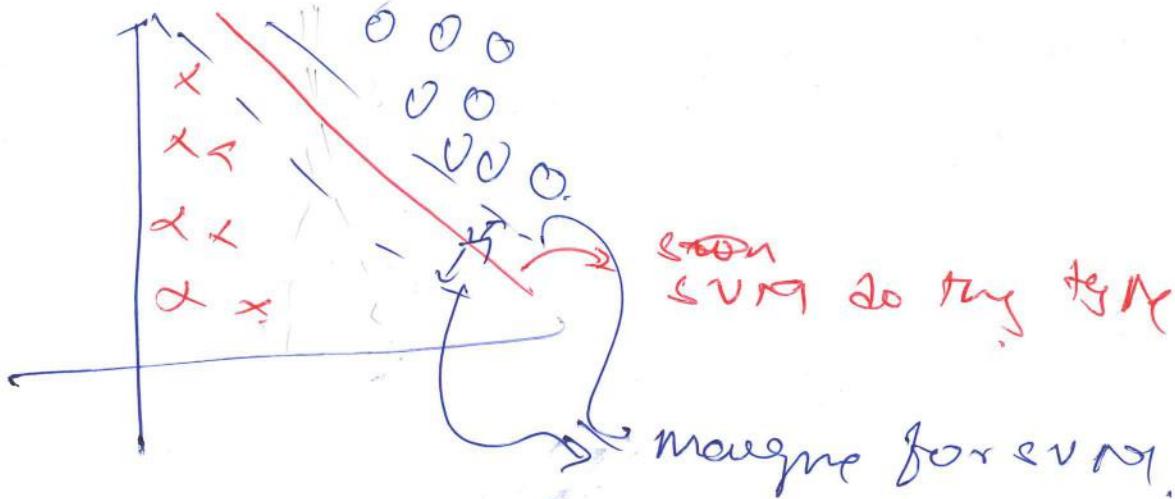
If $y=1$, we use $\theta^T x_i$ (not just ≥ 0), $\theta^T x_i \neq 0$
 If $y=0$, we want $\theta^T x_i = 1$ (not just ≥ 0) $\theta^T x_i = 1$

if $y=1$ does rule

SVM :- Boundaries

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

If $y^{(i)}=1$
 If $y^{(i)}=0$



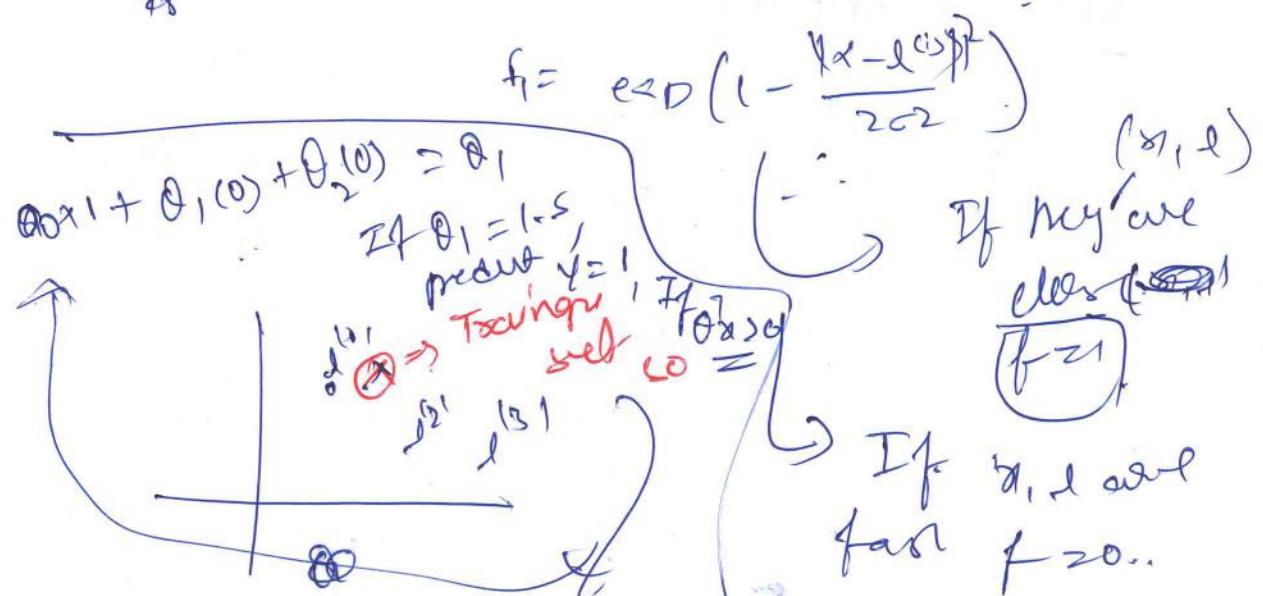
SVM with Kernels

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

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

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

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



$$\text{Training} \underset{\theta}{\min} C \sum_{i=1}^m g^{(i)} \text{cost}_i(\theta^T \phi^{(i)}) + (1-g^{(i)}) \text{cost}_0(\theta^T \phi^{(i)}) \\ + \frac{1}{2} \sum_{j=1}^{n+1} \theta_j^2$$

$$g(\theta_j) = \theta^T \phi$$

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

$$\rightarrow \theta^T \phi \cdot \theta \quad \left\{ \begin{array}{l} \text{making me equal} \\ \text{for reward.} \end{array} \right\}$$

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 fit very smooth, higher bias, low variance.

Logistic Regression vs SVM

n = number of features, ($\propto E2^{nH}$), 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:

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

→ If n is small, m is large

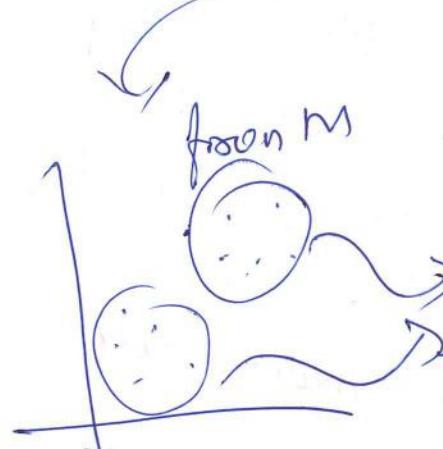
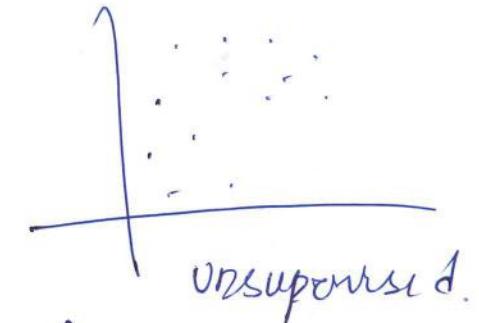
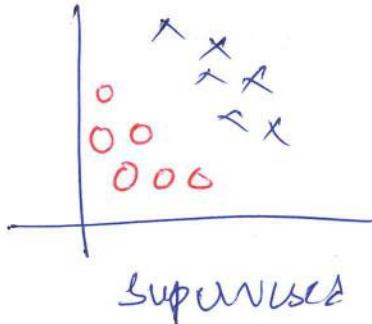
→ Overfit/odd not features, then use logistic regression & SVM without a kernel.

→ Neural networks works for most of these settings but may be slower to train.

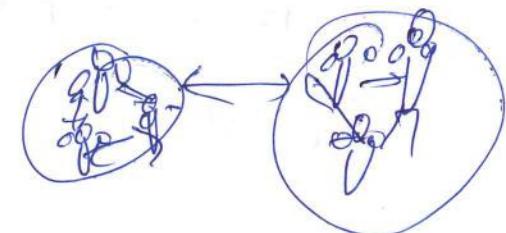
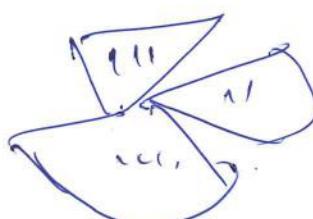
Unsupervised Learning

Unsupervised Learning (No labeling as in supervised)

Introduction :-



Assignment of clusters



Forward
propagation

- Clusters! K-means Algorithm.
-
- Each point is compared to cluster centroids
 - depends on them they are clustered
 - more iterations of k-means are done, to find K-number of clusters.

K-means algorithm.

Input:

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

$x^{(i)} \in \mathbb{R}^n$ (does $n=1$, something)

K-means algorithm

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

~~Test for convergence~~

repeat {

for $i = 1$ to m

$c^{(i)} =$ index (from 1 to K) 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 w.r.t.
 $u_k - u_k'$~~

y

~~repeat
until converge~~

~~update
centroids~~

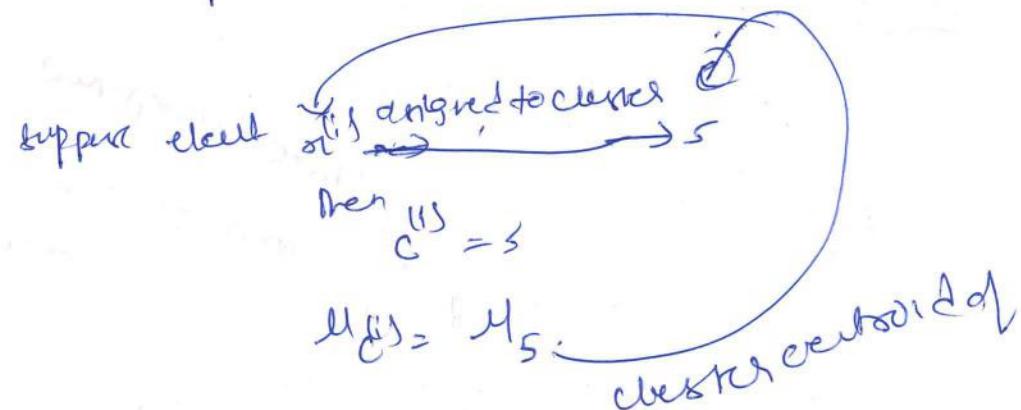
K-Means optimization Objective

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

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

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

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



Optimizatⁿ objective

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

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

distance between

centroid and object
has to be minimize & next

they can be almost clustered same.

K-Means Initialization

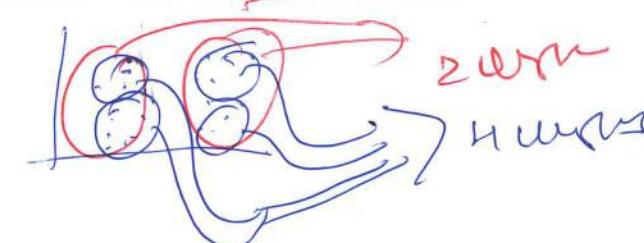
~~see driver for more~~

random initialization

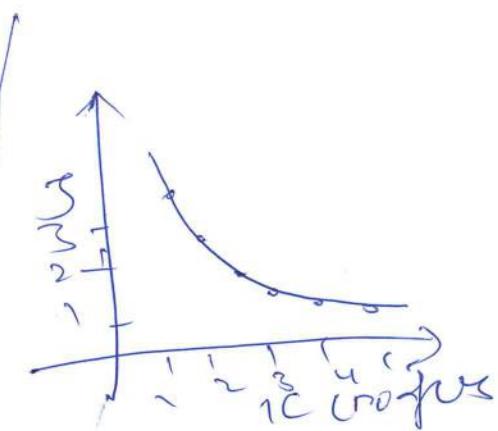
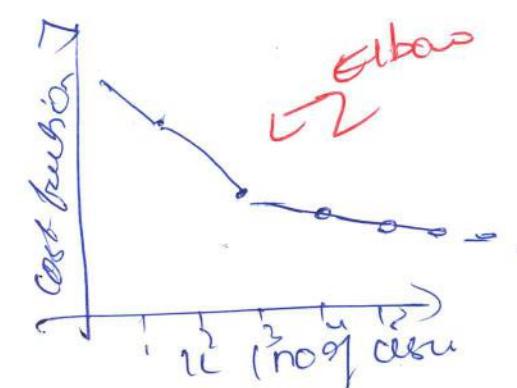
- should $n \leq m$.
- randomly pick ~~one~~ n training examples

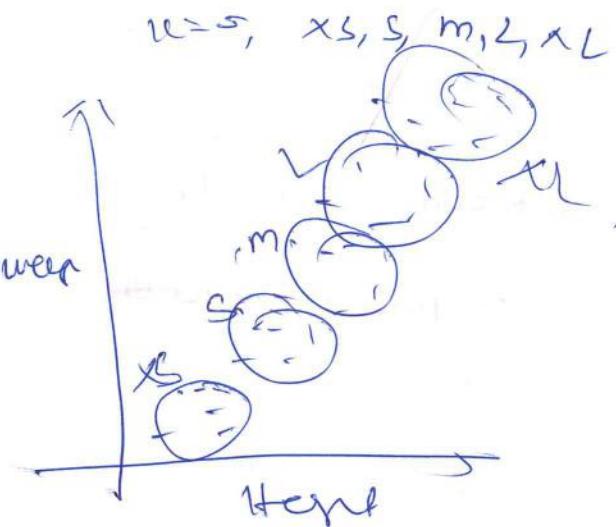
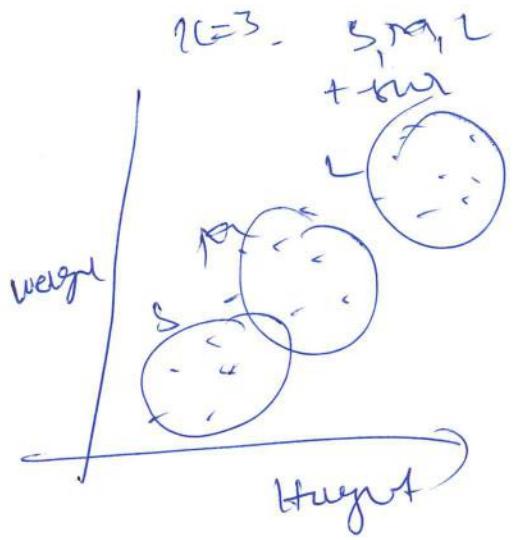
- set μ_1, \dots, μ_n equal to
these n samples

choose m number of cluster

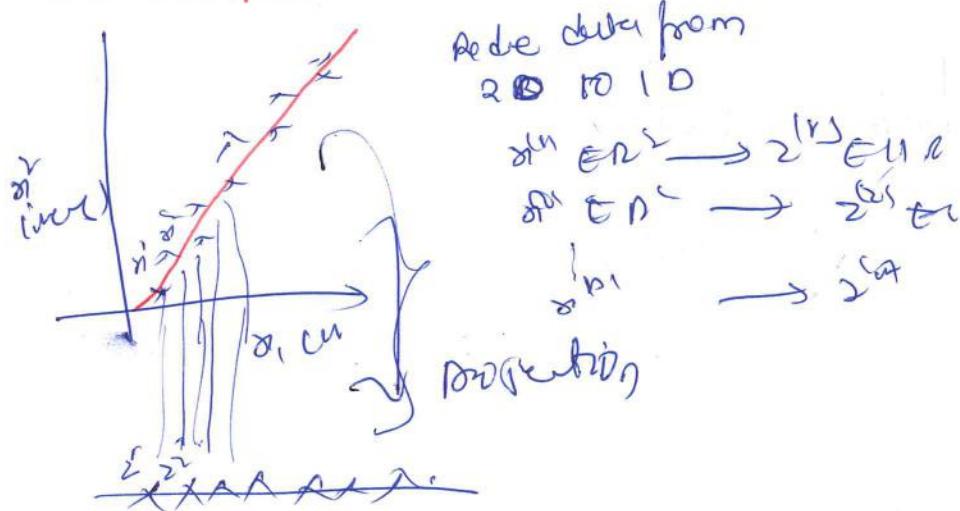


Elbow method



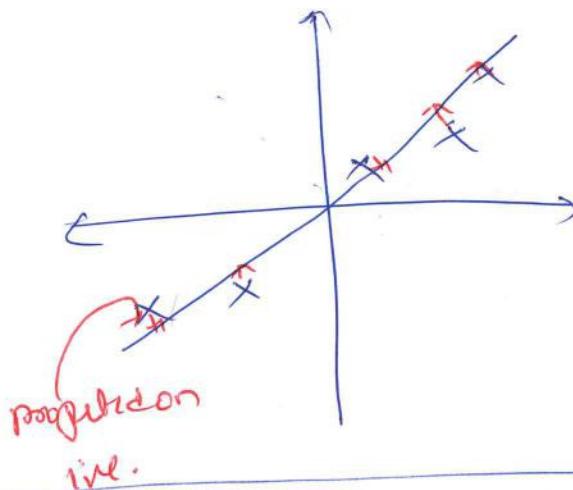


Data compression



PCA
DIMENSIONALITY REDUCTION:
PRINCIPAL COMPONENT ANALYSIS.

PCA: Formulation



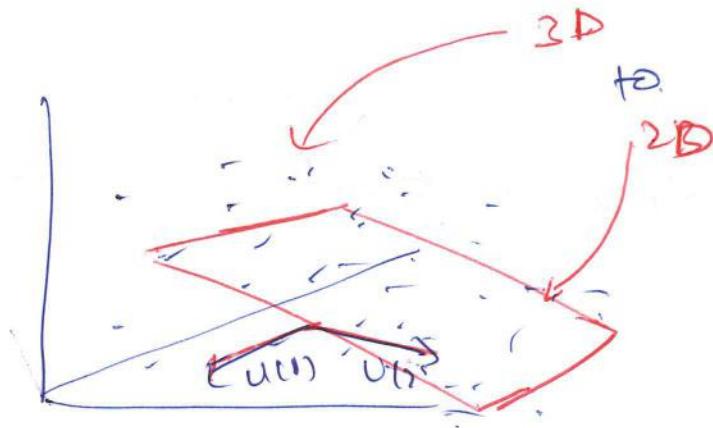
PLA :- Reduce 2D - to - 1D with min projection.
Find direction (a vector $u \in \mathbb{R}^n$).

General case

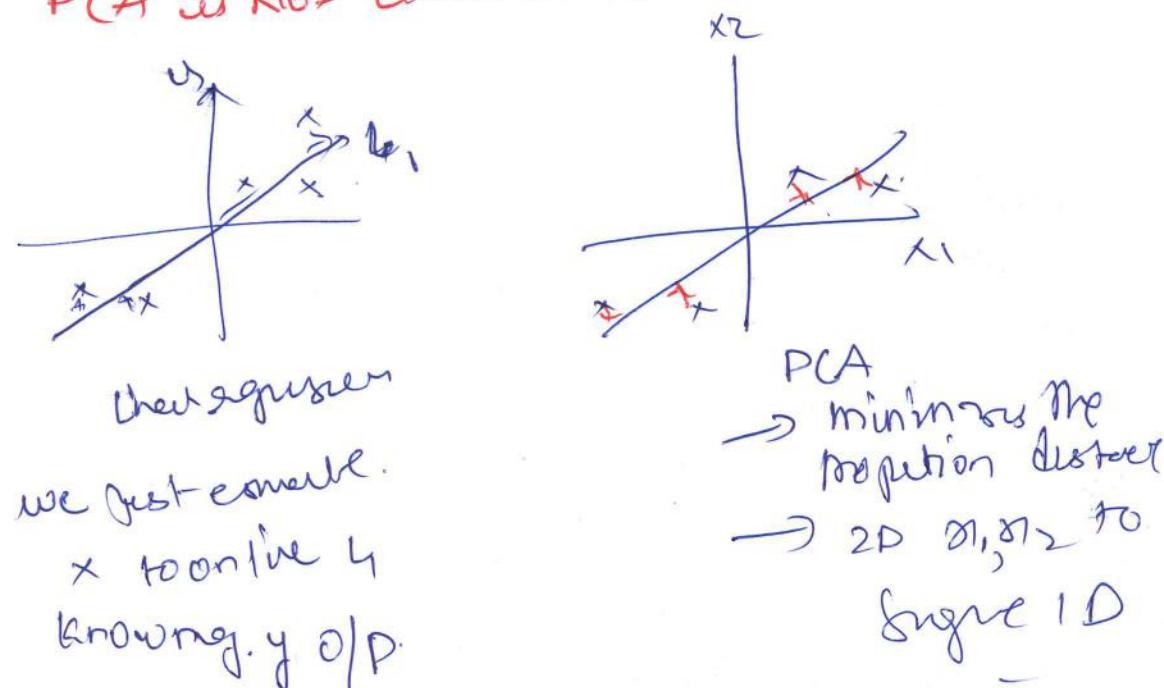
Reduce n -dimension to k -dimensions: find
 k - Vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project
the data, so as to minimize the projection error.

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

Q1



→ PCA is NOT linear regression



PCA Algorithm

Data Preparation of Data set: $\{x^{(1)}, \dots, x^{(m)}\}$
 scaling norm (mean Normalization)

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

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 data from n-D to 1-D
 compute "covariation 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

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

$U \cdot x \rightarrow 0$

~~take~~
 $x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^m$

$$z = \begin{bmatrix} 1 & 1 & 1 \\ u_{11} & u_{12} & \dots & u_{1n} \end{bmatrix}$$

Mathematical Proof Beyond the Scope
nxn
U ordered

$$(U_{\text{order}})_{x^k}^T = \begin{bmatrix} -N^{(k)} \\ -V_{1,2,\dots} \\ \vdots \\ -V_{m,n} \end{bmatrix} \quad (m)$$

Summary of PCA

Algorithms

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

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

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

$$z = U_{\text{order}}^T * x$$

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

$$\hat{x}_{\text{app}} = \frac{1}{m} \sum x^{(i)}$$

Reconstruction from compressed representation!

$$z = U_{\text{order}}^T * x$$

$$x = \underset{\text{approx}}{U_{\text{order}}} \cdot z$$

on = nxn real

choosing k (Node PCA)

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

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

* 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% Variance ~~greater~~ retained

Algorithm:- (PCA Charts)

Try PCA with $k=1$

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

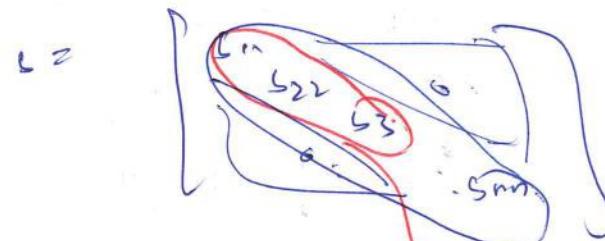
$x_i^{(1)} \text{ approx} \rightarrow x_i^{(m)}$

Check if

$$\frac{1}{m} \sum_{i=1}^m \|(\bar{x}_i^{(1)} - \bar{x}_i^{(m)})\|^2 \leq 0.01$$

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

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



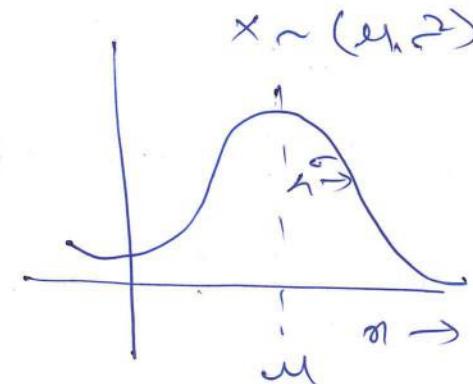
for que 1c

$$1 - \frac{\sum_{i=1}^K S_{ii}}{\sum_{i=1}^m \|x_i\|^2} \leq 0.01$$

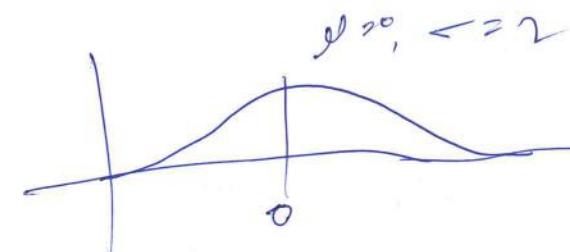
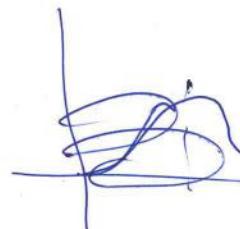
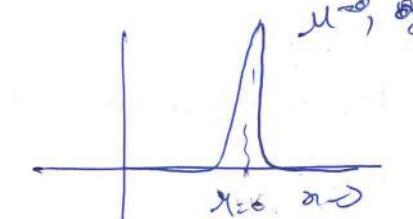
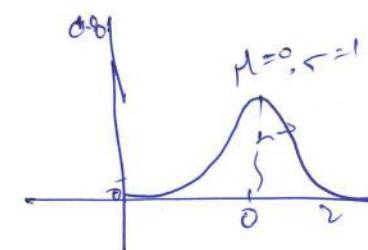
Density Estimation:

Gaussian distribution

$n \in \mathbb{R}$, 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 : $\{x_1^{(1)}, \dots, x_n^{(m)}\}$

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

(x_i, u_i, c_i)

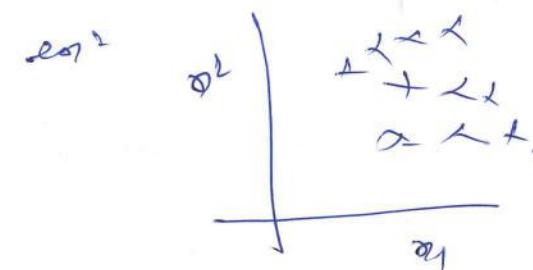
$$P(\sigma) = P(x_1, u_1, c_1) \cdot P(x_2, u_2, c_2) \cdots P(x_n, u_n, c_n)$$

↙ ↗ ↘

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 fault is out-of-trend)

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 $P(x|D)$

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

PREDICTIVE MOVIE RATINGS:

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

n_u = no of users

Problem Formulation :-

$n =$

$\rightarrow s(i,j) = 1$, if user i has rated movie j
(0,0 means)

$\rightarrow y_{(i,j)} = \text{rating by user } j \text{ for movie } i$
 (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: $\hat{r}_{j(i)} = \theta_0 + \sum_{k=1}^m \theta_k s_{j(k)}$

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

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

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

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

$y_i = b_i - \theta^T x_i$

	θ_0	θ_1	θ_2	θ_3	θ_4	θ_5	θ_6	θ_7	θ_8	θ_9	θ_{10}
Dom 1	5	2	3	1	2	3	2	1	2	3	2
Dom 2	2	1	2	3	2	1	2	3	2	1	2
Dom 3	0	1	0	2	1	0	2	1	0	2	1
out	0	0	0	0	0	0	0	0	0	0	0

(?) → finding max for user
perturbed movie rating

Optimization problems

Optimization problems

To learn $\theta^{(j)}$ (Parameters for $v_{0(j)}$):

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

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

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

so finally:

optimization algorithm:

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

$\mathcal{J}(\theta_1, \dots, \theta^{(n)})$

gradient descent update:

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

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

$$[00] \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \end{bmatrix} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

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

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

$$\theta_1 = 1.5 \\ \theta_0 = \frac{1.5}{3} = 0.5$$

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

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

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

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

→ Minimize both simultaneously.

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

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

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

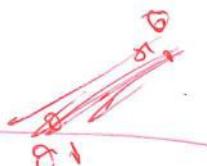
Collaborative Filter algorithm

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

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

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

$$\frac{d}{d \theta_{1j}^{(i)}} J$$



Low Rank Matrix Factorization

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

predicted rating

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

$$\begin{array}{c} (\theta^{(1)})^T (x^{(1)}) \quad \theta^{(2)} (x^{(2)}) \quad \dots \quad (\theta^{(n)})^T (x^{(n)}) \\ \theta^{(1)} (x^{(1)}) \\ \vdots \\ (\theta^{(n)})^T (x^{(n)}) \end{array}$$

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

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

SUMMARY OF COURSE

- Andrew Ng

Main topics

$$\left[\begin{array}{c|c} \theta^{(1)T} & x^{(1)} \\ \vdots & \vdots \\ \theta^{(n)T} & x^{(n)} \end{array} \right] = X \Theta^T$$

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.

ML Applying tools

④ Advice on building a machine learning system

~~Bias/variance, regularization; deciding what to work on next~~

evolution of learning algorithm, debugging, error analysis, testing analysis.

172
Court
Campbell