

Machine Learning

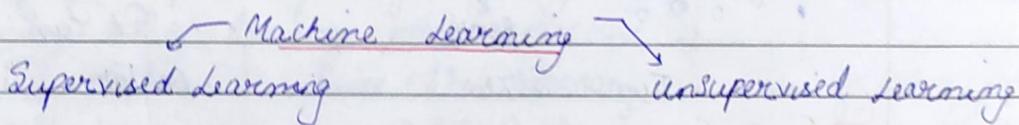
- By Andrew Ng

- * Tom Mitchell (1998) - A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E .

Ex: Checkers : $E \rightarrow$ Playing 1000's of games
 $T \rightarrow$ Play checkers
 $P \rightarrow$ winning or losing Probability.

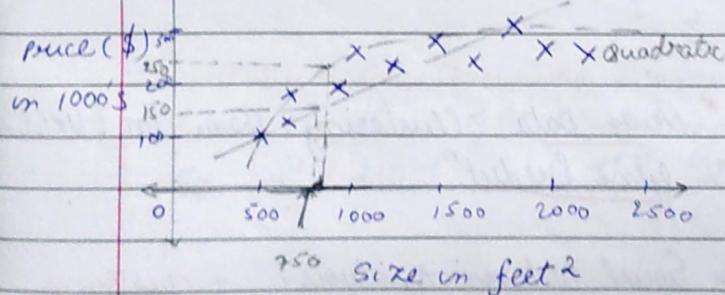
Ex: Spam & Non-spam Email - ID *

Classifying = Task T
 Labelling emails = Experience E
 Correctly classified = Performance P



- * Supervised learning = "Right answers" given to the host

Ex: Housing price prediction = regression - predict continuous valued output (price)



Ex Breast Cancer (^{Harmful}_{Malignant}, ^{Harmless}_{benign}) :

Malignant

1(Y)

0(N)

x x x x x

Classification :

discrete valued

output (0 or 1)

Tumor Size

There can be more classification too like

There are 3 types of cancer
 Benign type 1 type 2 type 3 cancer

o o o x o x o x x

Tumor Size

Ex

Age

This will be
 under benign side
 benign side
 o o x o
 o o o o
 o o o o

Malignant

Benign

Tumor Size

We can also use different parameters like

- Clump thickness

- Uniformity of cell size

- Uniformity of cell shape

... It can use ∞ no. of features too

Ex of Supervised Learning

* Questⁿ of Classification & Regression

Hacking \rightarrow output Yes or No \rightarrow Classification

Say produce crops which will sell or not \rightarrow Regression

* Unsupervised Learning :

(0 0)
 (0 0 0)
 (0 0 0)
 (0 0)
 (0 0 0)
 (0 0 0)
 (0 0)

Clustering

* New Data : Clustering them in 1 yell link

* Gene Predict

Ex Social Network Analysis -

Market Segments

Astronomical Data Analysis,

Organizing Computer Cluster

{ clustering problems }

~~1-10 in English~~

Non-clustering

Speaker 1 Nick 1 \leftarrow diff Postⁿ
 Speaker 2 \sim 1-10 in french Nick } cocktail party problem

To differentiate the two voice

We use this single line of code

[W, S, V] = svd((repmat(sum(x.*x, 1), size(x, 1), 1) .* x) .* x');

* cocktail party prob Identify individual voice

- * We will be using Octave as our environment
- * question for unsupervised learning algorithm
 clustering like → News articles group them on same story
 → Database for discovering market segments.

Intro Quiz 1 x 2 ✓ Regresⁿ 3 x Axisⁿ 4 x 5 ✓ definitⁿ Regresⁿ

Model and Cost function:

We saw that Example of housing System

Training set of housing prices (Portland)

Size in feet ² (x)	Price (\$) in 1000's (y)	m = No. of training ex.
2104	460	x's = input variable / features
1416	232	
1534	315	y's = output variable / target variable
852	178	
...	...	
(x	y	

(x, y) = one training example

(x⁽ⁱ⁾, y⁽ⁱ⁾) = ith training example

ex x⁽¹⁾ = 2104, y⁽²⁾ = 232.

Training sets



Learning Algorithm

Size of house $\rightarrow h \rightarrow$ Estimated price predicted

x hypothesis y
 h maps from x 's to y 's

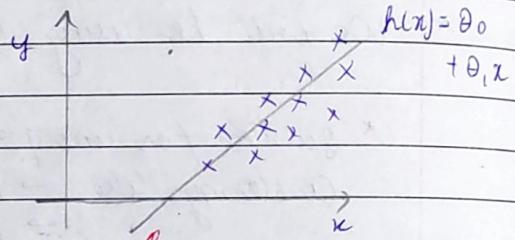
* How do we represent h ?

$$h_0(x) = \theta_0 + \theta_1 x$$

Short hand

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

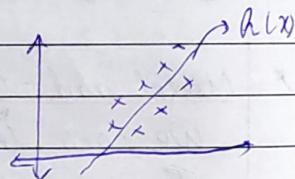
Linear Regression with one variable. }
 Univariate Linear Regression }



Name of this model.

⇒ Using the data set in previous page

$$\text{Hypothesis} \rightarrow h_0(x) = \theta_0 + \theta_1 x$$

 θ_i 's = ParametersHow to choose θ_0, θ_1 so that $h_0(x)$ is close to y for our training example (x, y)

$$\text{minimize}_{\theta_0, \theta_1} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

↑ training set
 ↓ outcome provided
 ↓ $\theta_0 + \theta_1 x^{(i)}$
 Mathematical convenience

$$\text{Cost funct} \rightarrow J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 \quad \left\{ \begin{array}{l} \text{Squared Error} \\ \text{funct} \end{array} \right.$$

$$\text{Goal} \rightarrow \text{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

$\underbrace{\text{cost f}}$

$$\text{parameters} \rightarrow \theta_0, \theta_1$$

\Rightarrow Simplified Hypothesis Equation:

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \text{or} \quad \theta_0 = 0$$

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

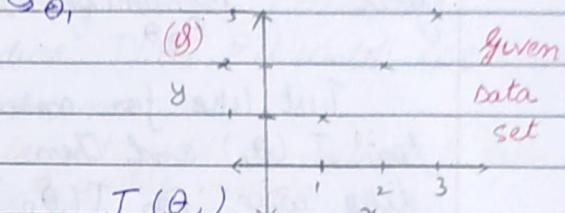
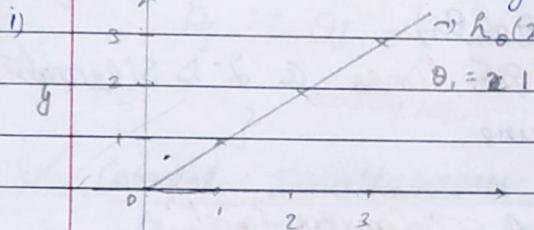
Parameters $\rightarrow \theta_0, \theta_1$

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

Goal \rightarrow minimize $J(\theta_0, \theta_1)$

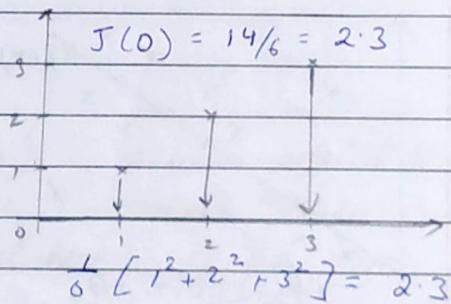
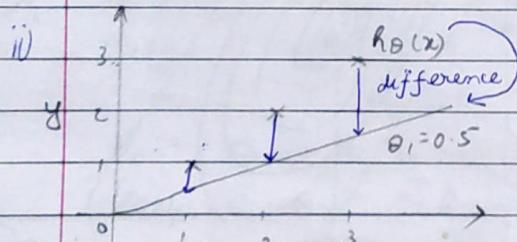
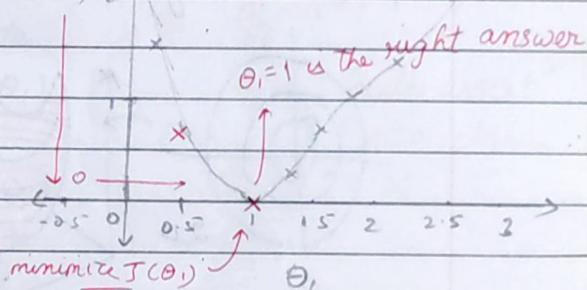
$$h_{\theta}(x)$$

for fixed θ_1 , this is a function of x function of the parameter θ_1 .



$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_{\theta_1}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0^2 \end{aligned}$$

$$\therefore \underline{J(1) = 0}$$



$$\underline{J(0.5)} = \frac{1}{2m} \left[\left(\frac{1}{2} - 1 \right)^2 + (1 - 2)^2 + \left(\frac{3}{2} - 3 \right)^2 \right]$$

$$= \frac{1}{2 \times 3} \left[\frac{14}{4} \right] = 0.58$$

$$\underline{J(0.5)} = 0.58$$

Now =

Hypothesis $h_0(x) = \theta_0 + \theta_1 x$

Parameters θ_0, θ_1

Cost f^n $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$

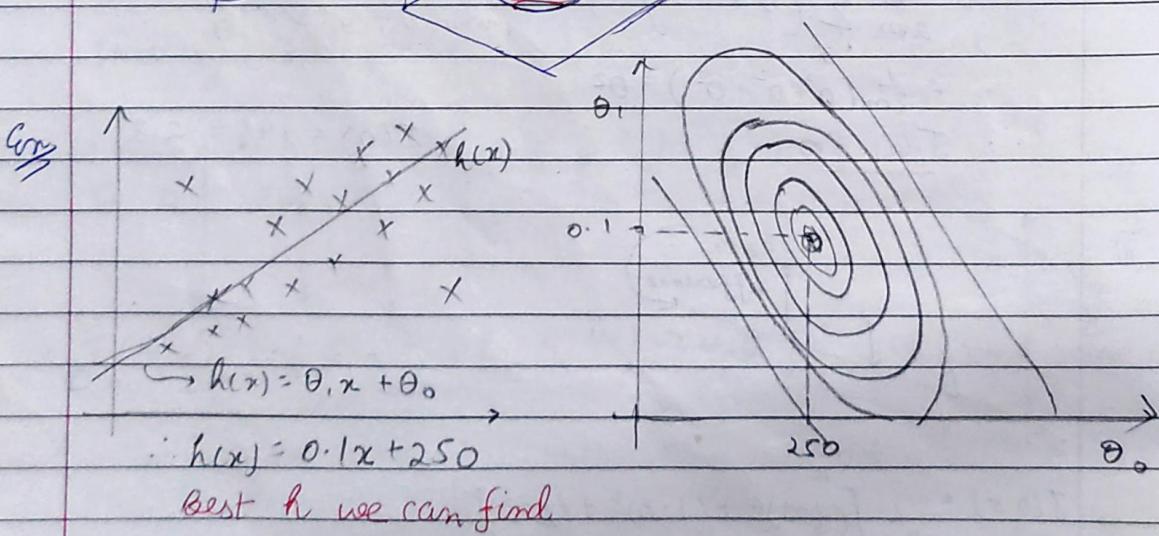
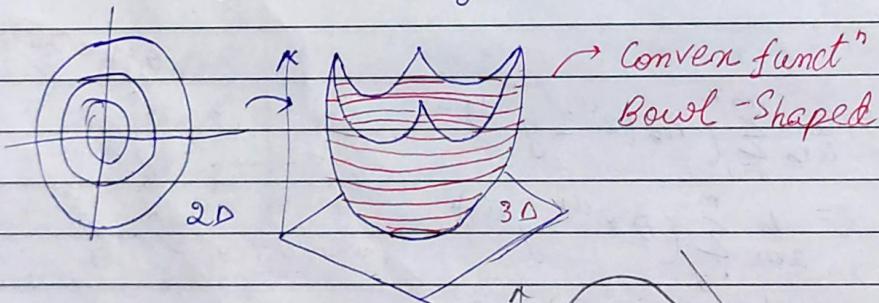
goal minimize $J(\theta_0, \theta_1)$

Just like for one variable θ , we had a 2-D dig for $J(\theta_0)$ and then we found global minimum.

Like wise for $J(\theta_0, \theta_1)$ we get a 3-D shape

~~Ex~~ $J(\theta_0, \theta_1) = z = x^2 + y^2$ we do the same analysis & we find the $\min J(\theta_0, \theta_1)$

We use for Contour lines a 2-D schematic representation of a 3-D figure



* Gradient Descent :-

We have some $f'' J(\theta_0, \theta_1)$ for $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

We want $\underset{\theta_0, \theta_1}{\min} J(\theta_0, \theta_1)$ or $\underset{\theta_0, \theta_1, \theta_2, \dots, \theta_n}{\min} J(\theta_0, \theta_1, \dots, \theta_n)$

⇒ Outline:-

Start with some θ_0, θ_1 [say $\theta_1 = \theta_0 = 0$]

keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum.

Gradient Descent Note its method to take a small steps towards a global minima.

repeat until convergence {

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

} Learning rate $\frac{\partial}{\partial \theta_j}$ Simultaneously update $\theta_0 \& \theta_1$ $\alpha := b$ Assignment operator

$$a := a + 1$$

✓ Correct : Simultaneous update

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

$a = b$ } Truth assert
Just a claim

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

$$X a = a + 1 X$$

$$\begin{aligned} \theta_0 &:= \text{temp}0 \quad \{ \text{Update } \theta_0 \& \theta_1 \\ \theta_1 &:= \text{temp}1 \end{aligned}$$

X Incorrect :-

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

then updated θ_0 will move to the further

$$\theta_0 := \text{temp}0 \quad \text{If we write it here}$$

steps that we don't want .

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

so make sure to update it at end .

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

Q) Suppose $\theta_0 = 1$, $\theta_1 = 2$ & we update it using the rule
 $\theta_j := \theta_j + \sqrt{\theta_0 \theta_1}$ (for $j=0 \& j=1$) what are the resulting values of $\theta_0 \& \theta_1$?

$$\Rightarrow \therefore \theta_0 := \theta_0 + \sqrt{\theta_0 \theta_1}$$

$$\theta_0 := \theta_0 + \sqrt{2}$$

$$\theta_0 := 1 + \sqrt{2} = 1 + \sqrt{2}$$

$$\theta_0 := 1 + \sqrt{2}$$

$$\theta_1 := 2 + \sqrt{2} < 2 + \sqrt{2}$$

$$\theta_1 := 2 + \sqrt{2}$$

$$\therefore \theta_0 = 1 + \sqrt{2} \& \theta_1 = 2 + \sqrt{2}$$

* gradient descent algorithm:

repeat until convergence δ

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

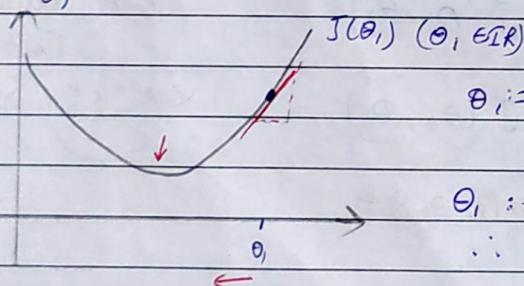
β

learning rate

$\frac{\partial \theta_j}{\partial \theta_j}$ derivative

$$\min_{\theta_1} J(\theta_1)$$

for $\theta_1 \in \mathbb{R}$ [real.]

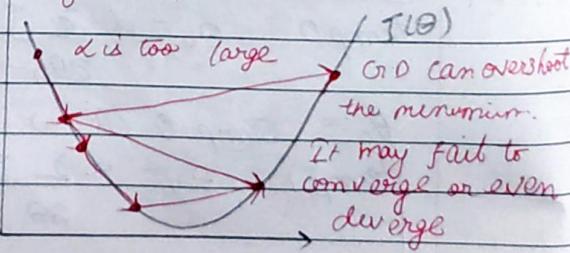
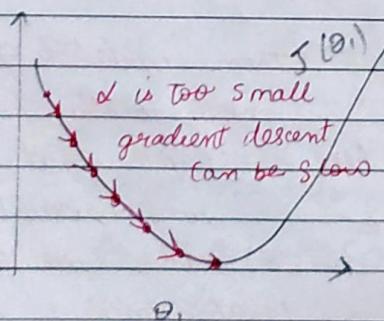


$$\theta_1 := \theta_1 - \alpha \frac{\partial J(\theta_1)}{\partial \theta_1} \quad \text{slope is (+)}$$

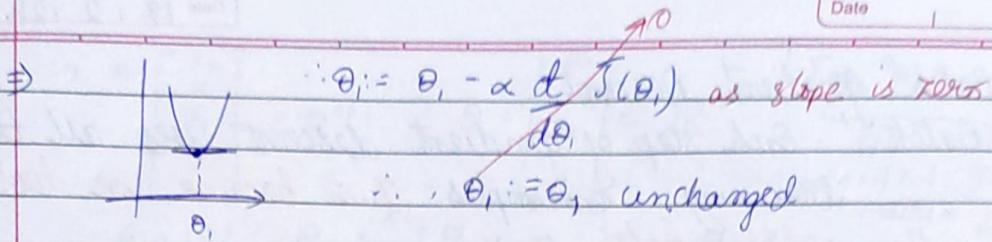
$$\theta_1 := \theta_1 - \alpha \text{ (positive number)}$$

\therefore we have to go towards left on scale which is true.

Same for opposite side:



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

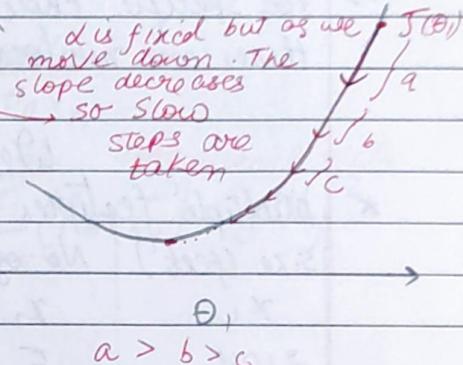


\Rightarrow Gradient Descent can converge to a local minimum, even with the learning rate α fixed.

$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$

\downarrow keeps on decreasing
 \downarrow α is fixed but as we move down, the slope decreases so slow steps are taken

As we approach a local minimum, $J(\theta)$ will automatically take smaller steps. So no need to decrease α over time



Gradient Descent Algorithm

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

for $j = 1$ & $j = 0$

Linear Regression Model

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

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

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

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

$$= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$\theta_0; j=0: \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

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

Gradient Descent Algorithm:

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

update θ_0 & θ_1 simultaneously

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

* Gradient Example

"Batch" gradient Descent

"Batch": Each step of gradient descent uses all the training examples. Just because we look at the whole batch (m)

- ⇒ Gradient descent can converge even if α is kept fixed
(But α cannot be too large, or else it may fail to converge)
- ⇒ for specific choice of cost function $J(\theta_0, \theta_1)$ in linear regression there are no local optima (other than global optimum)

Week - 2

* Multiple features (variables) :-

Size (feet ²)	No. of bed	No. of floors	Age of home years	Price (\$1000)
x_1	x_2	x_3	x_4	y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

$n = 4$

Notation :-

n = number of features

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

$x_j^{(i)}$ = Value of feature j in i^{th} training example

$$\rightarrow x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$$\rightarrow x_3^{(2)} = 2$$

Hypothesis :- New

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$\text{Ex: } h_{\theta}(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$$

example \uparrow base price \uparrow beds \uparrow price \uparrow price \downarrow price
 \uparrow price \uparrow floor \uparrow size \uparrow age of home

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

for convenience of notation; define $x_0 = 1$ $[x_0^{(i)} = 1]$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

Take dot and
we get matrix $h_{\theta}(x)$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$\underline{h_{\theta}(x) = \theta^T x} \rightarrow \text{Multi-variate linear regression}$$

* Gradient descent for Multiple variables :

$$h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

Parameters = $\theta_0, \theta_1, \dots, \theta_n / \theta_{n+1}$ - dimensional vector

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

Gradient descent

repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) \quad \left. \begin{array}{l} \text{(simultaneously update} \\ \text{for every } j=0, \dots, n \end{array} \right\}$$

gradient descent

Preciously ($n=1$):

Repeat {

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

$$\left. \begin{array}{l} \frac{\partial}{\partial \theta_0} J(\theta) \sum_{i=1}^m x_0 = 1 \end{array} \right\}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

}

New Algorithm ($n \geq 1$):

Repeat {

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

$$\left. \begin{array}{l} \text{(simultaneously update } \theta_j \text{ for} \\ j=0, 1, 2, \dots, n \end{array} \right\}$$

$$\text{like } \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

* Gradient Descent in Practice

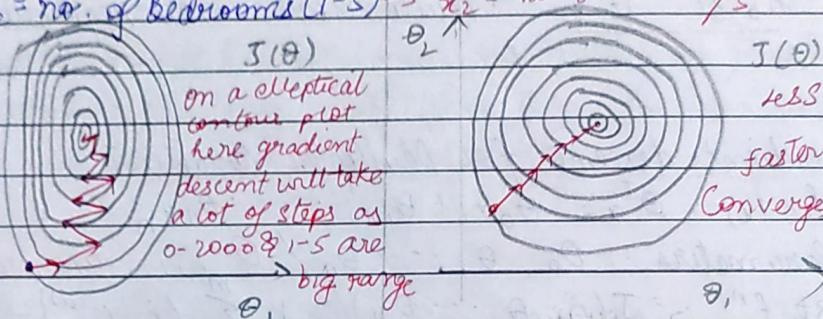
1. Feature Scalling:

Idea: Make sure features are on a similar scale.

$$\text{Ex } x_1 = \text{size (0-2000 feet}^2) \rightarrow x_1 = \frac{\text{size (feet}^2)}{2000}$$

$$x_2 = \text{no. of bedrooms (1-5)} \rightarrow x_2 = \text{no. of bedrooms}/5$$

We must do
feature scaling
to make it work
so that less steps
are taken.



Ignore θ_0

Feature scaling: Get every feature into approximately

$$x_0 = 1 \quad -1 \leq x_i \leq 1 \text{ range}$$

good $\{0 \leq x_i \leq 3\}$ $\{x_i - 100 \leq x_i \leq 100\}$ poor scaling
 scaling $\{-2 \leq x_2 \leq 0.5\}$ $\{x_4 - 0.001 \leq x_4 < 0.001\}$ Not approx to
 will work $\{-3 \leq x_i \leq 3\}$ $\{-\frac{1}{3} \leq x_i \leq \frac{1}{3}\}$ works

⇒ mean normalization:

Replace x_i with $x_i - \mu_i$ to make features have approx. zero mean [Do not apply to $x_0 = 1$]

$$(0-2000) \text{ Ex } x_1 = \frac{\text{size} - 1000}{2000} \rightarrow -0.5 \leq x_1 \leq 0.5$$

$$(1-5) \quad x_2 = \frac{\# \text{bedroom} - 2.5}{5} \rightarrow -0.4 \leq x_2 \leq 0.4$$

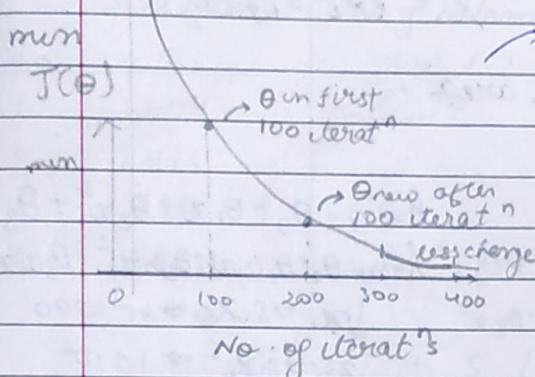
$$\therefore x_1 \rightarrow \frac{x_1 - \mu_1}{s} \quad \begin{matrix} \text{avg value of} \\ x_1 \text{ in set} \end{matrix}$$

$$(3) \rightarrow \text{range (max-min or standard deviation)}$$

2. Learning Rate =

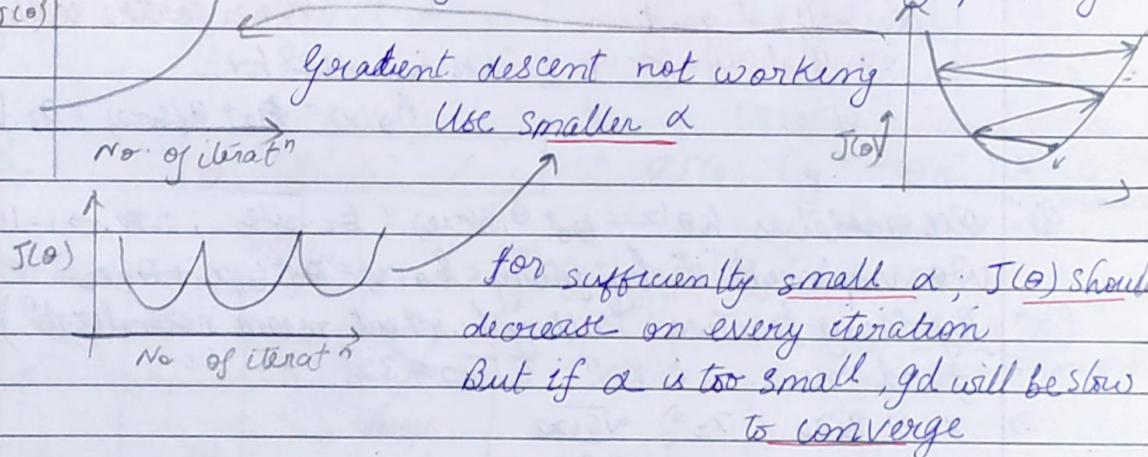
$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}, \text{ repeat } m \text{ of iterations}$$

Debugging = How to make sure gradient descent is working correctly
How to choose learning rate α



Example automatic convergence-test
declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration

\Rightarrow cases to make sure gradient descent is working properly:-



\Rightarrow Summary:-

If α is too small \Rightarrow slow convergence

If α is too large \Rightarrow $J(\theta)$ may not decrease on every iteration;
may not converge

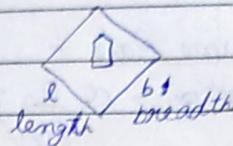
Try $\alpha = 0.001, 0.01, 0.1, 1, \dots$
grad = $\frac{1}{m} \sum_{i=1}^m h(\theta^{(i)}) - y^{(i)} \cdot x^{(i)}$

* Features & polynomial Regression :

Ex Housing price prediction

So our $h_{\theta}(x)$ will be

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



So instead using $l \& b$ we can substitute them with Area

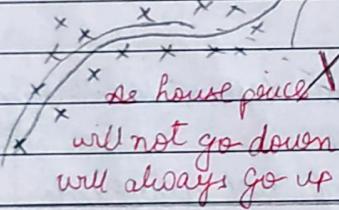
& make a new variable $X = l \times b = \text{Area}$

$$h_{\theta}(x) = \theta_0 + \theta_1 X$$

↓ (and area)

* Polynomial Regression :

price
(y)



will not go down
will always go up

size (Z)

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$h_{\theta}(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

$$x_1 \rightarrow \text{size} \Rightarrow 1 - 1000$$

$$x_2 \rightarrow \text{size}^2 \Rightarrow 1 - 10^6$$

$$x_3 \rightarrow \text{size}^3 \Rightarrow 1 - 10^9$$

Hence feature scaling becomes imp.

We can also $h_{\theta} =$

$$h_{\theta}(x) = \theta_0 + \theta_1 / \text{size} + \theta_2 \sqrt{\text{size}}$$

Q) Our model is $h_{\theta}(x) = \theta_0 + \theta_1 / \text{size} + \theta_2 \sqrt{\text{size}}$; size $\rightarrow 1 - 1000$ (feet)²

We implement it by fitting $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

By doing feature scaling (without mean normalization)
find new $x_1 \& x_2$ ($\sqrt{1000} \approx 32$)

$$\Rightarrow x_1 = \frac{x}{1000}; x_2 = \frac{\sqrt{x}}{32}$$

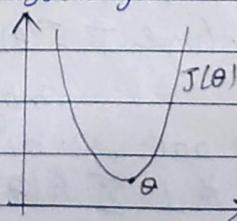
* Normal Equation : Method to solve for θ analytically

Intuition : If $\nabla J(\theta) = 0$

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

$$\text{So we do } \frac{d}{d\theta} J(\theta) = 0$$

& solve for θ



But we need to solve it for

$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } j) \quad \& \text{ solve for } \theta_0, \theta_1, \theta_m$$

Size (feet) ²	No. of bedrooms	No. of floors	Age of home (years)	Price 1000\$
x_0	x_1	x_2	x_3	y
2104	5		45	460
1416	3		40	232
1534	3	2	30	315
852	2	1	36	178

$$m=4$$

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$$M \times (n+1)$$

$$\theta = (X^T X)^{-1} X^T y \rightarrow \nabla J(\theta) = (y^T - \theta X^T) X = 0$$

$$y^T X - \theta X^T X = 0$$

m-example $\hat{y} = (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ - n features $X^T X = \theta X^T X$

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad X = \begin{bmatrix} \cdots (x^{(1)})^T \\ \cdots (x^{(2)})^T \\ \vdots \\ \cdots (x^{(m)})^T \end{bmatrix} \quad \theta = y^T X \underbrace{(X^T X)^{-1}}_{\text{Pseudo inverse}}$$

Ans If $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix} \therefore X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

$$\therefore \text{Ans } \theta = (X^T X)^{-1} X^T y$$

\Rightarrow format long | format short | $A = [1 \ 2 \ ; \ 3 \ 4 \ ; \ 5 \ 6]$

a:

$a = 3.14159265358979$ | $a = 3.1416$ | $A = 1 \ 2$
| $3 \ 4$
| $5 \ 6$

$\Rightarrow v = [1 \ 2 \ 3] \quad z = [1; 2; 3] \quad v = 1:0.2:2$
 $v = 1 \ 2 \ 3 \quad z = \frac{1}{3} \quad v = 1 \ 0 \ 1 \cdot 2 \ 1 \cdot 4 \ 1 \cdot 6 \ 1 \cdot 8 \ 2 \cdot 0$

$\Rightarrow v = 1:6 \quad ones(2,3) \quad c = 2 \ ones(2,3) \quad w = zeros(1,3)$
 $v = 1 \ 2 \ 3 \ 4 \ 5 \ 6 \quad ans = 1 \ 1 \ 1 \quad c = 2 \ 2 \ 2 \quad w = 0 \ 0 \ 0$

$\Rightarrow rand(3,3) \quad randn(1,3) \quad w = 6 + sqrt(10) * (randn(1,1000))$
 $ans = 0.8 \ 0.3 \ 0.7 \quad w = -1.4 \ -1.2 \ -0.6 \quad hist(w)$
 $0.1 \ 0.2 \ 0.3 \quad (Creates a histogram)$
 $0.9 \ 0.8 \ 0.4 \quad his(w, 50)$



$\Rightarrow eye(4) \quad help eye \quad help rand \quad A = [1 \ 2 \ ; \ 3 \ 4 \ ; \ 5 \ 6];$

Diagonal matrix | brings up documentation | size(A)

$\begin{matrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{matrix}$

for both the commands

Help Help

ans 3 2

sz = size(A);

size(sz)

ans 1 2

$\Rightarrow v = [1 \ 2 \ 3 \ 4]^{*4}; \quad pwd \quad cd (\text{to change})$

length(v)
ans = 4

(location/path)
(Mathlib Dirive) directory

ls (gives list of direct-
ories on my desktop)

$\Rightarrow load features.dat \quad who \quad size(features) \quad who \rightarrow$ gives

load('features.dat') Variable in the ans = 47 2 details of all the

to load any file current workspace matrix mxn of that data variable

$\Rightarrow clear features \quad v = price(1:3) \quad save hello.mat v; \quad clear$

removes the 5 looks 1st, 3 elements | save v in hello.mat | it clears
loaded data 12 of some loaded folder on desktop | everything
data price

\Rightarrow Save hello.txt v -ascii; $A = [1\ 2\ 3\ 4\ 5\ 6]$; $A(:, 2)$
 Save data as text
 $\{ A(2, :) \rightarrow$ means every element along
 $\{ \text{ans} = 3\ 4$ that rows / col. 0

get from row 3 new value
 \Rightarrow ans = 1 2 5 6 | $A = \begin{bmatrix} 1 & 10 \\ 2 & 11 \\ 5 & 12 \end{bmatrix}$ assign new value to the matrix | $A = [A, 100]$ append another column
 vector to right
 \Rightarrow $A = [1; 2; 5; 6] \quad C = [A; B]$ size(C)
 ans = 1 2 3 4 5 6 | $B = [4; 5; 6]$ 2 5 ans = 1 on top ans = 6 |
 int as single vector | $C = [A \ B]$ 3 6 4 of each |
 concatenates A B B 5 other | 6 other

$$\Rightarrow A = [1 \ 2; 3 \ 4; 5 \ 6] \quad | \quad A * C \quad | \quad A * B \text{ multiply} \quad | \quad x \cdot 2 \\ B = [11 \ 12; 13 \ 14; 15 \ 16] \quad | \quad \text{multiply} \quad | \quad \begin{matrix} 11 & 12 \\ 13 & 14 \\ 15 & 16 \end{matrix} \text{ individually} \quad | \quad \text{ans } 1 \ 4 \quad \text{Square} \\ C = [1 \ 1; 2 \ 2] \quad | \quad A * C \quad | \quad \begin{matrix} 39 \\ 75 \end{matrix} \text{ with} \quad | \quad 9 \ 16 \quad 16 \\ \text{values of matrix} \quad | \quad 25 \ 36 \quad \text{matrix}$$

$v = [1; 2; 3]$	$\log(v)$	$\exp(v)$	$\text{abs}(v)$	$-v$
1 / \sqrt{v} <u>real part</u>	0.00	2.71	1	cm^{-1}
ans 1:00	0.693	7.38	2	-2
0.500	1.098	20.08	3	-3
0.250				

$\Rightarrow V + \text{ones}(\text{length}(V), 1)'$ $V + 1$ A' Transpose (A') matrix
 ans 2. & add 1 in OR ans 2. $a_1 / 3 5$ $a_2 / 4 5$
 $\begin{matrix} 3 \\ 4 \end{matrix}$ all the numbers $\begin{matrix} 3 \\ 4 \end{matrix} 2 4 6 \quad \begin{matrix} 3 \\ 4 \end{matrix} 6$

$\Rightarrow a = [1 \ 15 \ 2]; \ [val, ind] = \max(a); \ a < 3 \ \text{find}(a < 3); \ \text{magic}(3)$

$Val = \max(a), Val = 15$

~~ans~~ $1 \ 0 \ 1$ ~~ans~~ $1 \ 3$, $8 \ 1 \ 6$ \leftarrow

~~ans~~ $Val = 15, ind = 2$ $T/F \uparrow$ $3 \ 5 \ 7$

$\uparrow \ 9 \ 2$ magic square

$\Rightarrow \text{magic}(A);$ $a = [1 \ 15 \ 2]; \max(a, [1, 1]) \sum(a, 1)$

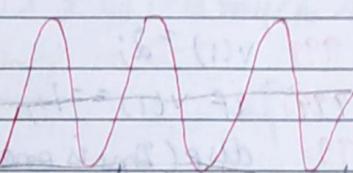
$[n, c] = \text{find}(A >= 7)$ sum(a) 18 ans 897 15 15

ans n = 1 c = 1 re(1,1)
 $\begin{matrix} 1 \\ 3 \\ 2 \end{matrix}$ $\begin{matrix} 1 \\ 2 \\ 3 \end{matrix}$ $\begin{matrix} 1 \\ 3 \\ 2 \end{matrix}$ prod(a) 30 max(a, [1], 2) sum(A, 2) 15
 $\begin{matrix} 2 \\ 3 \end{matrix}$ $\begin{matrix} 2 \\ 3 \end{matrix}$ floor(a) 1/15 2 ans 8/9 col wise 15

$\Rightarrow \text{eye}(3);$ sum(sum(A * eye(9))) flip up! eye(3), inv V(A) invert

$A * \text{eye}(3)$ ans = 15 0 0 1 flip 0.147 -0.14 0.06
 $\begin{matrix} 8 & 0 & 0 \\ 0 & 8 & 0 \\ 0 & 0 & 2 \end{matrix}$ 0 1 0 it -0.06 0.02 0.10
 1 0 0 -0.01 0.18 -0.10

$\Rightarrow x = [0: 0.01: 1];$ 0.5
 $y_1 = \sin(2 * \pi * 4 * x);$ 0.1
 $\text{plot}(x, y_1)$ -0.1



$\Rightarrow y_2 = \cos(2 * \pi * 4 * x);$ hold on; to overlap new plot

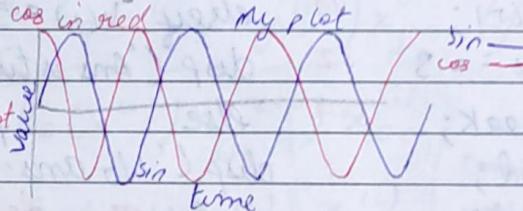
$\text{plot}(x, y_2, 'r');$

$x\text{label}('time')$

$y\text{label}('value')$

legend('sin', 'cos')

title('my plot')



$\Rightarrow \text{figure}(1); \text{plot}(x, y_1);$

$\text{figure}(2); \text{plot}(x, y_2);$

subplot(1, 2, 1);

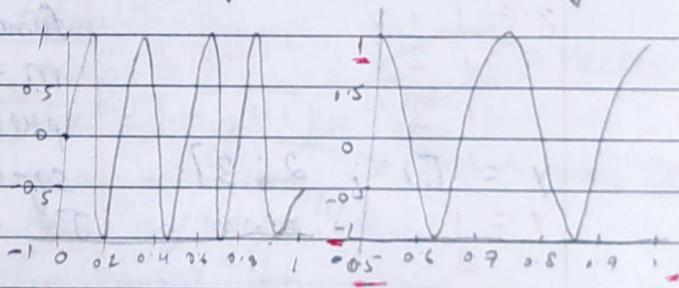
plot(x, y_1);

subplot(1, 2, 2);

plot(t, y_2);

axis([-0.5 1 -1 1])

x-axis y-axis



clf; clear figure

$\Rightarrow A = \text{magic}(5)$; * See the figure in Notes

$\text{magicSc}(A)$

$\text{imagesc}(A)$, colorbar, colormap gray;

$v = [2, 4, 8, 16, 32, 64]$

$\Rightarrow \text{for } i = 1:3, \text{ indices} = 1:3; \text{ for } i = \text{indices}, c = 1;$

$v(i) = 2^{\wedge} i;$	$i \in \text{indices}$	$(\text{disp}(i);$	$\text{while } i \leq 4$
$\text{end};$	$\text{index} = 1, 2, 3$	$\text{end};$	$v(i) = 99$
$* v = [2, 4, 8]$	$\{$	$\{$	$i = i + 1;$
			$\text{end};$

$\Rightarrow c = 1; \text{ ans } 999 \quad v(1) = 2;$

while True, $999 \quad \text{if } v(1) == 1,$

$v(c) = 999; \quad 99 \quad \text{disp('Ans is one');}$

$i = i + 1; \quad 99 \quad \text{else if } v(1) == 2,$

$\text{if } i == 3 \quad 32 \quad \text{disp('Ans is two');}$

$\text{break}; \quad 64 \quad \text{else}$

$\text{end}; \quad \text{clisp('No ans');}$

$\text{end}; \quad \text{Then Ans is two!}$

$\theta = [0; 1]$

\uparrow

$3, 3$

\downarrow

$2, 2$

\downarrow

$1, 1$

\downarrow

$0, 0$

\downarrow

$\theta = [0; 0]$

* Now we want to predict cost f^n using Octave:

$X = [1, 1; 1, 2; 1, 3]$

make a notepad.m

function J = costFunction(X, y, theta)

m = Size(X, 1);

predictions = X * theta;

sqrErrors = (predictions - y).²;

J = 1/(2 * m) * sum(sqrErrors);

$y = [1; 2; 3]$

$y =$

2

3

$\theta = [0; 1];$

$j = \text{costFunction}(X, y, \theta)$

$j = 0$

$\theta = [0; 0];$

$j = \text{costFunction}(X, y, \theta)$

$j = 2.33$

$$\% i.e. \frac{1^2 + 2^2 + 3^2}{2 \times 3} = 2.33$$

$\frac{1^2 + 2^2 + 3^2}{2 \times 3}$

As understanding starts from

1

1

* Vectorization :-

$$g_{\theta}(x) = \sum_{j=0}^n \theta_j x_j \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$$= \theta^T x$$

Unvectorized Implementation

$$\rightarrow \text{prediction} = 0.0;$$

$$\rightarrow \text{for } j = 1:n+1,$$

$$\rightarrow \text{prediction} = \text{prediction} +$$

$$\theta_j x_{(j)} \quad \text{is much more efficient.}$$

$\rightarrow \text{end};$ *vectorized imp in C++

Vectorized Implementation

$$\text{prediction} = \theta^T x;$$

both will give same

output But vectorized imp.

\Rightarrow Gradient Descent :-

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

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

$n=2$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m h_{\theta}(x^{(i)}) - y^{(i)} x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m h_{\theta}(x^{(i)}) - y^{(i)} x_2^{(i)}$$

, vector subtraction

simultaneous

updates

Vectorized Implementation :-

$$\text{Then} \rightarrow \theta := \theta^{TRn+1} - \alpha S^{TRn+1}$$

$$\text{where } S = \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x^{(i)}$$

Vector

real Number IR

$$S = \begin{bmatrix} S_0 \\ S_1 \\ S_2 \end{bmatrix} \quad S_0 = \frac{1}{m} \sum h_{\theta}(x^{(i)}) - y^{(i)} x_0^{(i)}$$

$$S_1 = \frac{1}{m} \sum h_{\theta}(x^{(i)}) - y^{(i)} x_1^{(i)}$$

what does summation says

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \end{bmatrix}$$

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

$$= (h_{\theta}(x^{(1)}) - y^{(1)}) \tilde{x}^{(1)} + (h_{\theta}(x^{(2)}) - y^{(2)}) \tilde{x}^{(2)} + \dots$$

Number vector

$$\text{like } u(j) = 2v(j) + 5w(j) \Rightarrow u = 2v + 5w \text{ (for all } j)$$

Week - 3

* Logistic Regression Classification :

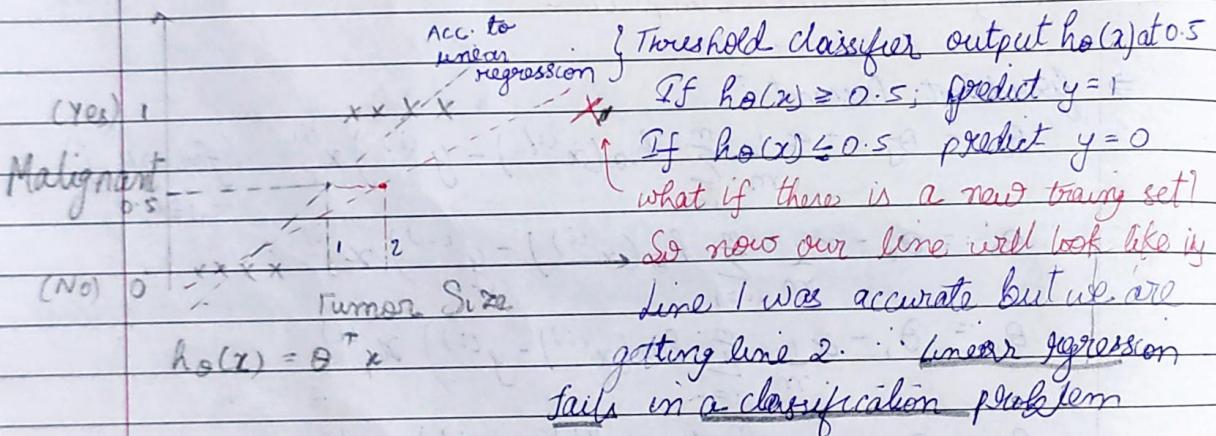
Ex Email: Spam / Not spam?

Online Transactions: Fraudulent (Yes/No)

Tumor: Malignant / Benign?

$y \in \{0, 1\}$ 0: "Negative class" (e.g., benign tumor)
1: "Positive class" (e.g., malignant tumor)

$y \in \{0, 1, 2, 3\} \rightarrow$ Multiclassification problem.



\Rightarrow Classification : $y = 0$ or 1

$h_0(x)$ can be > 1 or < 0

We will make a model

Logistic Regression : $0 \leq h_0(x) \leq 1$

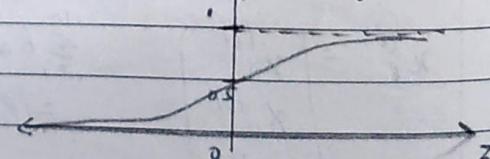
\hookrightarrow classification algorithm

In linear regression we had $h_0(x) = \theta^T x$

In logistic regression $h_0(x) = g(\theta^T x)$

$$g(x) = \frac{1}{1 + e^{-x}} \quad \left\{ \begin{array}{l} \text{sigmoid function} \\ \text{logistic function} \end{array} \right.$$

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



Parameters θ

⇒ Interpretation of hypothesis output

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

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

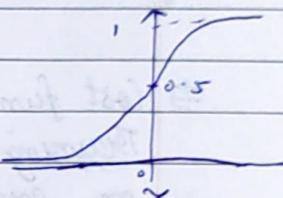
$$h_{\theta}(x) = 0.7 ; y=1$$

Tell patient that 70% chance of tumor being malignant.

$h_{\theta}(x) = P(y=1 | x; \theta) \Rightarrow$ "prob. that $y=1$, given x , parameterized by θ "

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

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



⇒ Logistic Regression:

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

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

$$g(x) \geq 0.5 \quad \text{when } x \geq 0$$

Predict $y=1$ if $h_{\theta}(x) \geq 0.5$ $\theta^T x \geq 0$

$$y=0 \text{ if } h_{\theta}(x) \leq 0.5 \quad \theta^T x < 0$$

$$\therefore h_{\theta}(x) = g(\theta^T x) \geq 0.5 \text{ whenever } \theta^T x \geq 0$$

⇒ Decision Boundary:

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \quad \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

Predict " $y=1$ " if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$$x_1 + x_2 \geq 3$$

$$x_1 + x_2 \leq 3 ; y = 0$$

Des'n Bound is only a property of Hypothesis not of dataset

\Rightarrow Non-linear Decision Boundaries :

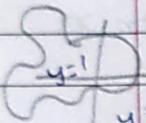
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad \therefore y = 1 \text{ if } -1 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 1 \text{ Circle}$$

We can often get more complex higher order features

$$\text{Like } h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



\Rightarrow Cost function for Logistic Regression :

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

m examples $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_0 = 1; y \in \{0, 1\}$

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

How to choose parameter θ ?

Cost functⁿ

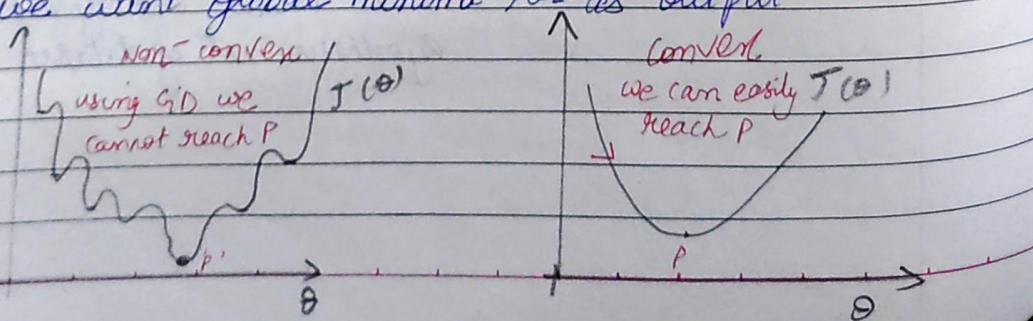
$$\text{linear regression } J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

We don't use square error because $(h_{\theta}(x) - y)^2$ and

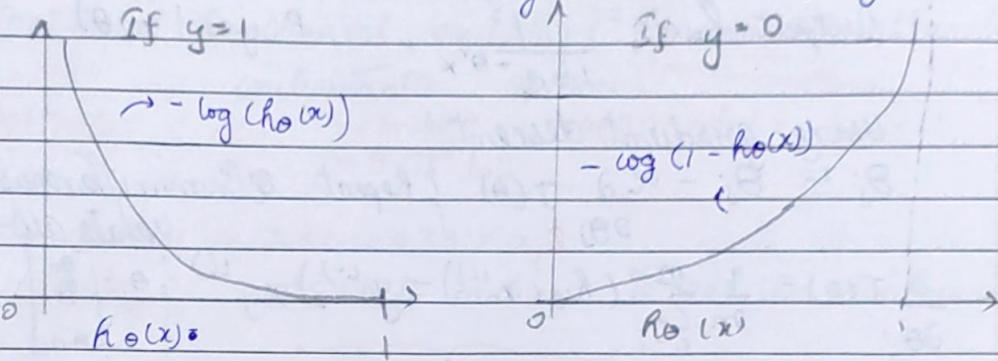
$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$ we won't get a good graph

we will have a graph with many local minima but
we want global minima for θ as output



⇒ Logistic Regression cost function :

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



cost = 0 ; if $y=1$ & $h_\theta(x) = 1$

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

If $h_\theta(x) = 0$ but $y = 1$

we'll penalize learning algo
by a very large cost.

⇒ Simplified cost function and gradient descent :

Logistic Regression cost function

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

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

$y=0$ or 1 always

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

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

$$\text{if } y=0 : \text{cost}(h_\theta(x), y) = -\log(1-h_\theta(x))$$

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

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

To fit parameters θ :

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

To make a prediction given new x :

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

Using gradient descent

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \text{Repeat } \quad \begin{matrix} \text{Simultaneously} \\ \text{update all } \theta_j \end{matrix}$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

Also looks identical to linear regression

* Advanced optimization for Logistic Regression :

We have other algorithms too except Gradient Descent.

- ⇒ Conjugate gradient
- ⇒ BFGS
- ⇒ L-BFGS

Advantages: No need to manually pick alpha
Disadvantage: Much slower than GD.

$$\text{Ex: } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad J(\theta) = (\theta_0 - 5)^2 + (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5) \quad ; \quad \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

We can clearly see $J(\theta)$ is min when $\theta_0, \theta_1, \theta_2 = 5$

We can solve this on Octave too :

function [J Val, gradient] = costFunction(theta)

Val = (theta(1)-5)^2 + (theta(2)-5)^2;

gradient = zeros(2, 1); (grad = [0 0])

gradient(1) = 2 * (theta(1)-5); diff w.r.t θ_1

gradient(2) = 2 * (theta(2)-5); \leftarrow

setting my parameters

- ⇒ Use this advance code :-

```

options = optimset ('GradObj', 'on', 'MaxIter', '100');
initialTheta = zeros (2,1); initializing θ; i.e. initial guess
[OptTheta, functionVal, exitFlag] = fminunc(@costFunction,
    initialTheta, options) fmin unconstraint → points to the
cost function
optTheta = 5.00      ↳ initial θ      ↳ options which we had set
5.00                ↳  $\theta \in \mathbb{R}^d, d \geq 2$ 

```

function Val = 1.577e-030 %

`exitFlag = 1 % help minimize %readThis documentation`

- ⇒ Basic framework for algorithm

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \rightarrow (\theta_0, \theta_1, \theta_2, \theta_3)$$

function [jVal, gradient] = costFunction(theta)

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

gradient(W) = [code to compute $\frac{\partial \text{loss}}{\partial W}(t)$];

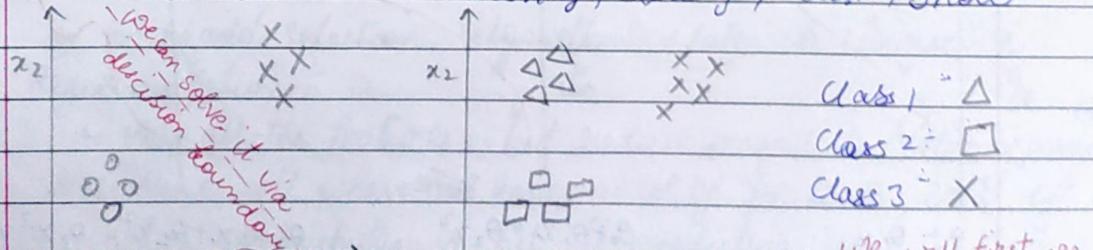
gradient(2) = [code to compute $\frac{\partial}{\partial \theta_j} J(\theta)$];

gradient(θ^{n+1}) = [code to compute $\frac{\partial}{\partial \theta} J(\theta)$];

- Multiclass Classification One vs all :

Ex Email foldering / tagging :- Work , Friends , Family ; Hobby
 $y = 1$ $y = 2$ $y = 3$ $y = 4$

Medical diagram :- Not ill, Cold, Flu
Weather :- Sunny, Cloudy, Rain, Snow

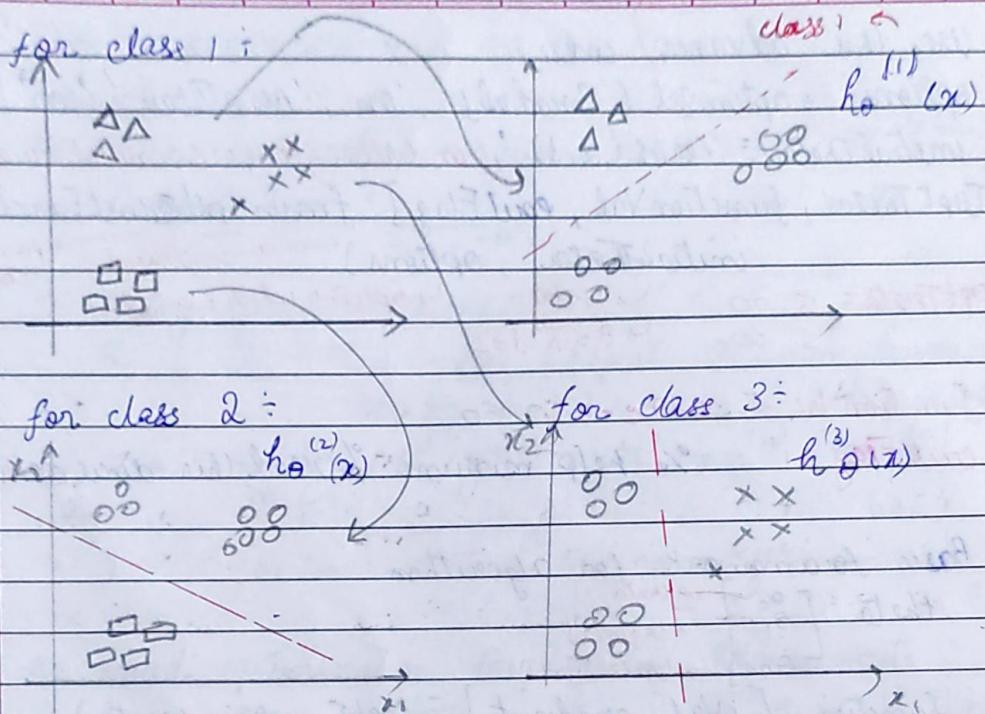


Binary Classification

Multi-class classification for class 1; rest

two we will make them same:

Q Treat as binary classification



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

Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class i to predict the probability that $y=i$.

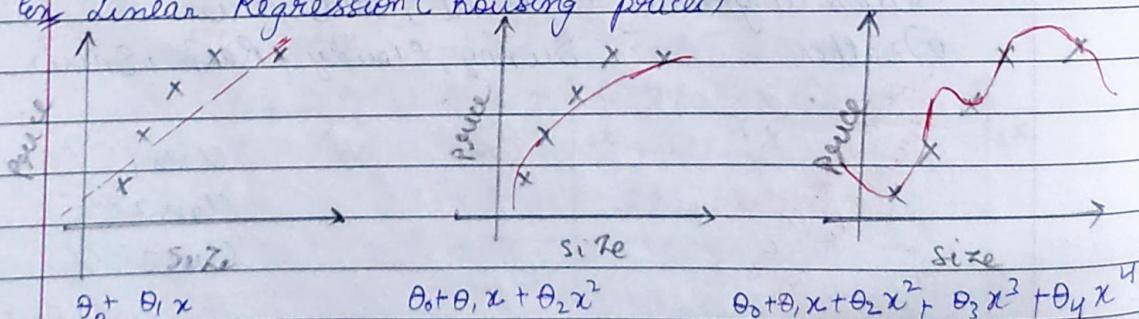
On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_\theta^{(i)}(x)$$

* Regularization :-

The problem of overfitting :-

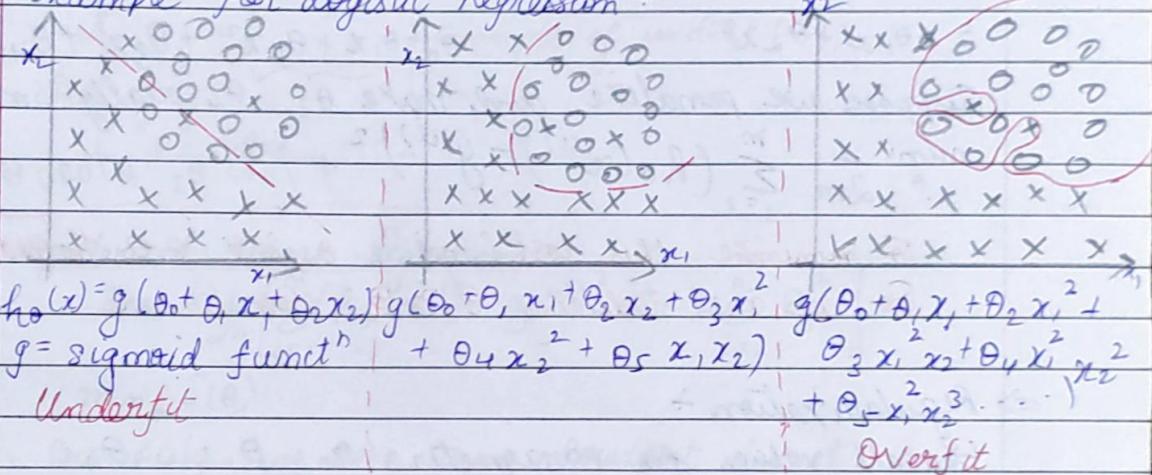
ex Linear Regression (Housing price)



Underfit "High Bias" "Just Right" Overfit "High Variance"

- ⇒ Overshooting : If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples)

- ⇒ Example for logistic Regression :



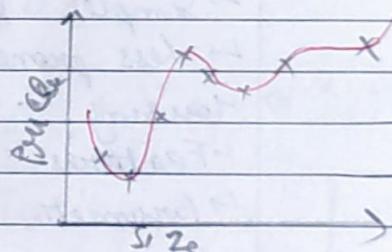
- ⇒ Addressing Overfitting :

x_1 = size of house x_4 = age of house

x_2 = No. of Bedrooms x_5 = kitchen size

x_3 = No. of floors

$x_{100} = \dots$



- ⇒ Options to resolve this overfitting :

- i) Reduce number of features :

↳ Manually select which features to keep

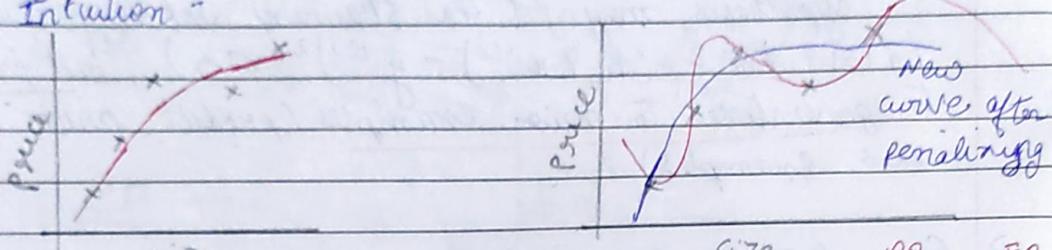
↳ Model selection algorithm (later in course).

- ii) Regularization :

↳ Keep all the features, but reduce magnitude / value of parameters

↳ Works well when we have a lot of features, each of which contributes a bit to predicting y .

\Rightarrow Regularization - cost function
Intuition -



$$y + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make θ_3, θ_4 really small

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{1}{1000} \theta_3^2 + \frac{1}{1000} \theta_4^2$$

To minimize this cost function $\theta_3 \approx \theta_4 \approx 0$ to reduce $J(\theta)$

So we are reducing $\theta_3 \approx \theta_4$ close to zero

\Rightarrow Regularization -

Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$

\hookrightarrow "Simpler" hypothesis $[\theta_3 \approx \theta_4 \approx 0]$

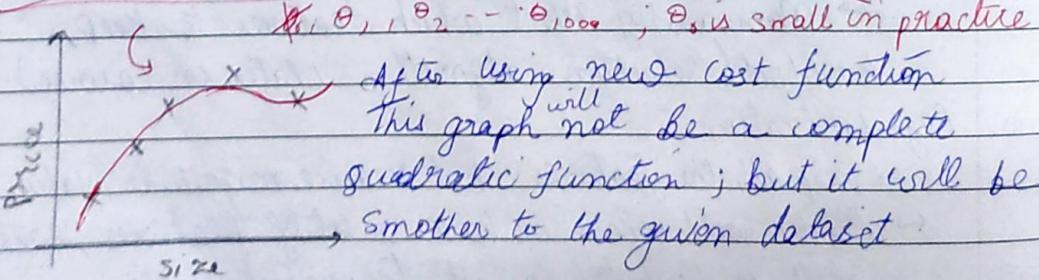
\hookrightarrow Less prone to overfitting

\Rightarrow Housing -

\hookrightarrow Features : x_1, x_2, \dots, x_{100}

\hookrightarrow Parameters : $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$ regularizⁿ parameter

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m h_{\theta}(x^{(i)}) - y^{(i)} \right]^2 + \lambda \sum_{j=1}^n \theta_j^2 \quad \begin{array}{l} \text{we don't take} \\ \theta_0 \text{ into consider} \end{array}$$



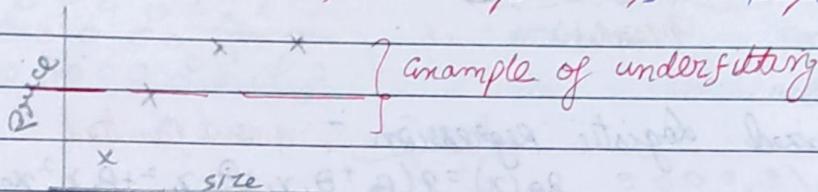
⇒ What if λ is set to an extremely large value (say $\lambda = 10^{10}$)?

If this happens then we have to eliminate θ_j

$$\text{Hence } \theta_1, \theta_2, \theta_3, \theta_4 \approx 0$$

So remaining hypothesis will be $h_\theta(x) = \theta_0$

$$h_\theta(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$



⇒ Regularized linear regression:

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

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

Previously we were using gradient descent

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

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

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

$1 - \alpha \frac{1}{m} < 1$
 m may be 0.99

⇒ Normal Equation

$$X = \begin{bmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix}_{m \times (n+1)} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}_{T R^m}$$

$$\min_{\theta} J(\theta); \frac{\partial J(\theta)}{\partial \theta_j} \stackrel{\text{set}}{=} 0$$

$$\theta = (X^\top X + \lambda \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & 0 \\ 0 & 0 & \cdots & 0 & 0 \end{bmatrix}_{(n+1) \times (n+1)})^{-1} X^\top y$$

Eg. $n=2$
 $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

\Rightarrow Non-invertibility \therefore

Suppose $m = n$,

#examples #features ; $\theta = (X^T X)^{-1} X^T y$

If $\lambda > 0$

$$\theta = (X^T X + \lambda \begin{pmatrix} 0 & 1 & & \\ & 0 & \ddots & \\ & & \ddots & 0 \\ & & & 0 & 1 \end{pmatrix})^{-1} X^T y$$

Invertible

\hookrightarrow non-invertible / singular
use pinv function

\Rightarrow Regularized logistic regression \therefore Sigmoid function

$$h(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_2^2 x_3 + \dots)$$

Cost function:

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

New cost function:

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

\Rightarrow Gradient Descent:

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

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

Not for $j = 0$

\Rightarrow Advance Optimization \therefore minimize (cost function)

function [jVal, gradient] = costFunction(theta)

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

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

gradient(1) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} + \frac{\lambda}{m} \theta_1$$

gradient(n) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];

Week - 4★ Neural Networks :

$$\begin{array}{ccccccc} x & x & x & x & x & x & x \\ \times & \times & x & x & x & x & x \\ x & x & x & x & x & x & x \\ x & x & x & x & x & x & x \\ x_2 & x & x & x & x & x & x \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \quad g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

when x_i is large

$x_1 = \# \text{size}$

$x_2 = \# \text{bedrooms}$

\vdots

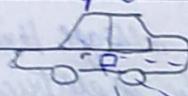
$x_{100} = \dots$

$i.e. x_1^2, x_1 x_2, x_1 x_3, \dots \approx O(n^2)$

But if we use features consisting of $x_1^2, x_2^2, x_3^2, x_4^2, \dots, x_{100}^2$ we will have 100 features so can't get shapes like ellipse. But we can never get this plot.

⇒ If we use cubic features $x_1 x_2 x_3, x_1^2 x_2, x_1 x_2 x_3, \dots$
 $O(n^3) \approx 170,000$ possible features.

⇒ Computer Vision :



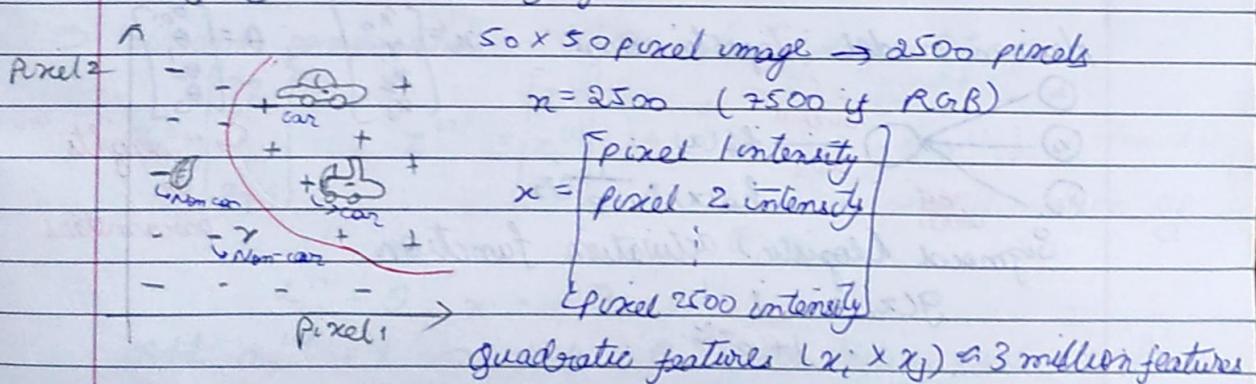
→ This is a car what humans see

we give our training set whether it is

1	5	9	4	8
7	2	1	3	
4	2	5	6	9

→ But the camera sees this
 specified that these are cars and these are not cars & then we give our testing set.

⇒ pixel → learning algorithm

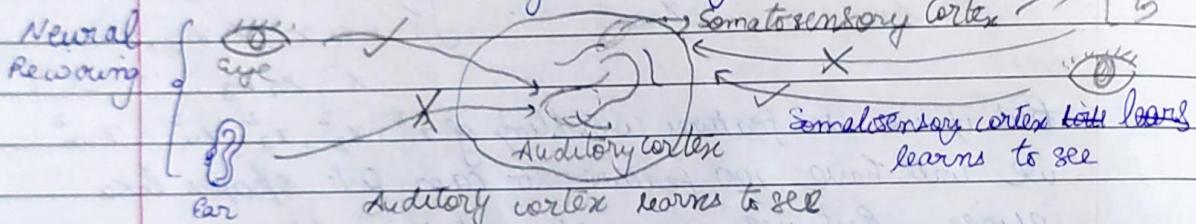


* Neurons and brains:

Neural Networks:

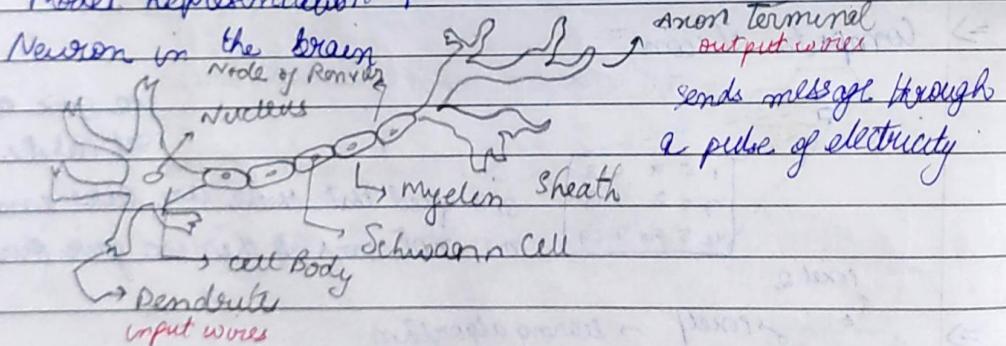
Origin → Algorithm that try to mimic the brain
 Very widely used in 80's and early 90's; popularity diminished in late 90's. Recent resurgence: State of the art technique for many applications.

⇒ The "one learning algorithm" hypothesis



** See notes for awesome examples of Neural Networks.

↗ Model Representation I:



neuron model: logistic unit

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

weights or parameters

$x_0 = 1$

x_1 bias unit

x_2 input wires

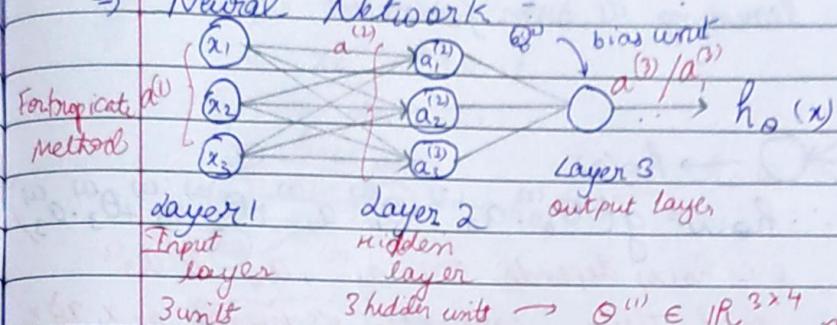
x_3 output wire

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

Sigmoid (logistic) activation function

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

⇒ Neural Network :



We can add $x_0/a_0^{(1)}$
not that important

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

$\Theta^{(j)}$ = matrix of weights controlling function mapping
hidden layer values from layer j to $j+1$

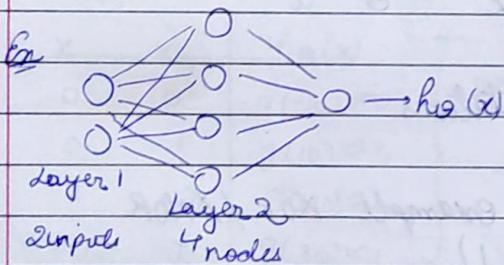
$$a_1^{(2)} = g(x_1^{(2)}) \quad a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3) \rightarrow x_1^{(2)}$$

$$a_2^{(2)} = g(x_2^{(2)}) \quad a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3) \rightarrow x_2^{(2)}$$

$$a_3^{(2)} = g(x_3^{(2)}) \quad a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3) \rightarrow x_3^{(2)}$$

$$h_0(x) = a^3 = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}) \rightarrow x^{(3)}$$

If network has s_j units in layer j ; s_{j+1} units in layer $j+1$; then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$



∴ find dimensions of $\Theta^{(1)}$

$$\cdot s_j = 2$$

$$\cdot s_{j+1} = 4$$

$$\cdot \text{Dimensions} = (s_{j+1}) \times (s_j + 1)$$

$$\Theta^{(1)} = 4 \times 3 \text{ dimensions}$$

for bias node a_0

⇒ Forward propagation = Vectorized Implementation:

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}; \quad x^{(2)} = \begin{bmatrix} x_1^{(2)} \\ x_2^{(2)} \\ x_3^{(2)} \end{bmatrix} \quad \begin{aligned} \alpha x^{(2)} &= \Theta^{(1)} x \\ \text{3dim vector } \alpha^{(2)} &= g(x^{(2)}) \in \mathbb{R}^{3 \text{dim. Vect}} \end{aligned}$$

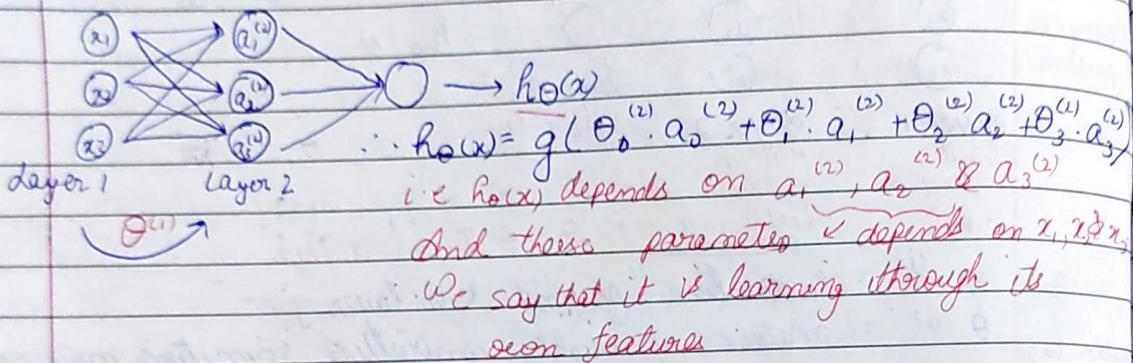
as $x^{(2)}$ is a product matrix of Θ & x .

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

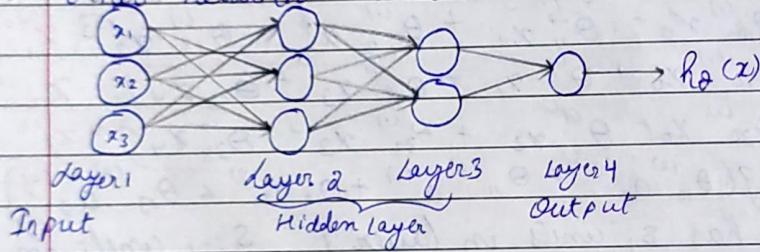
Add $a_0^{(2)} = 1$ \curvearrowright note: $\alpha^{(2)}$ is a 4dim feature vector

$$x^{(3)} = \Theta^{(2)} \alpha^{(2)}; \quad h_0(x) = a^{(3)} = g(x^{(3)})$$

⇒ Neural Network learning its own features :



⇒ Other network architectures :



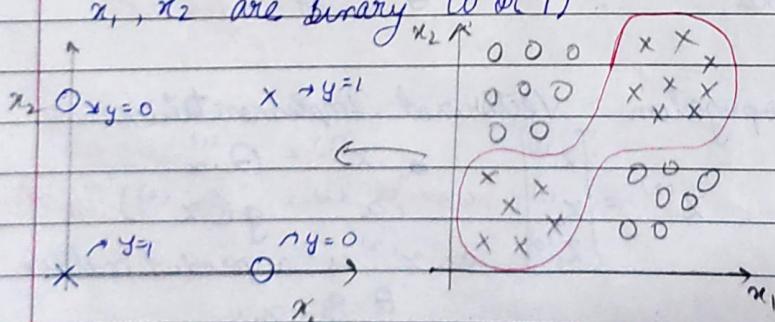
∴ How would you compute $a^{(2)}$?
 $\therefore a^{(2)} = g(x^{(2)})$ & $x^{(2)} = \theta^{(1)} a^{(1)}$

⇒ Neural Network : Representations :

Examples and Intuition !

Non-linear classification example : XOR / XNOR

x_1, x_2 are binary (0 or 1).



$$y = \neg x_1 \text{ XOR } x_2$$

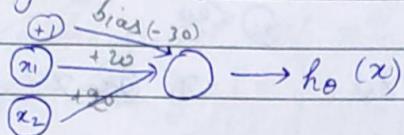
$$\hookrightarrow x_1 \text{ XNOR } x_2$$

i.e. NOT ($x_1 \text{ XOR } x_2$)

\Rightarrow Simple example : AND

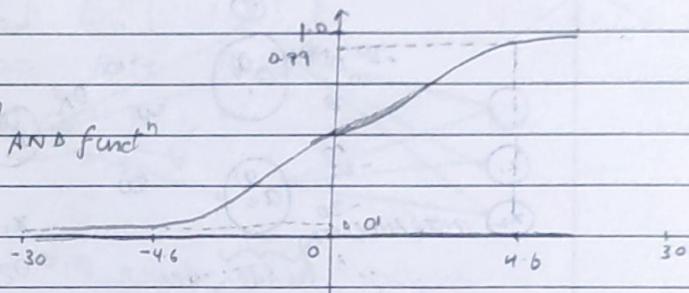
$$x_1, x_2 \in \{0, 1\}$$

$$y = x_1 \text{ AND } x_2$$



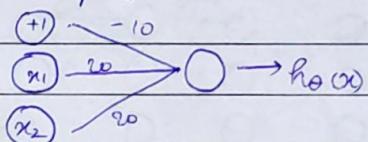
$$\therefore h_\theta(x) = g(-30 + 20x_1 + 20x_2)$$

x_1	x_2	$h_\theta(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$



$$\therefore h_\theta(x) \approx x_1 \text{ AND } x_2$$

\Rightarrow Example : OR

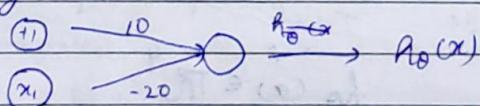


$$h_\theta(x) = g(-10 + 20x_1 + 20x_2)$$

x_1	x_2	$h_\theta(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$g(10) \approx 1$
1	1	$g(30) \approx 1$

$$h_\theta(x) \approx x_1 \text{ OR } x_2$$

\Rightarrow Negation = NOT



x_1	$h_\theta(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

$$h_\theta(x) = g(10 - 20x_1)$$

$$(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2) = 1$$

$$\text{If and only if } x_1 = x_2 = 0$$

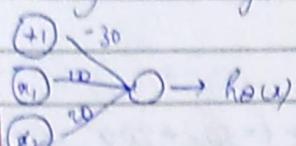
We can create complex
calculus $h_0(x)$ [Non linearities]
via Neural networks

Page No.

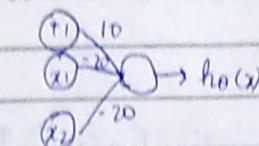
Date

80

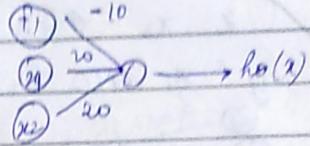
⇒ Putting it together: $x_1 \text{ XNOR } x_2$



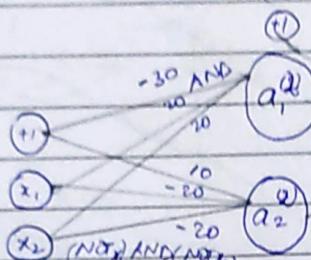
$x_1 \text{ AND } x_2$



(NOT x_1) AND (NOT x_2)



$x_1 \text{ OR } x_2$

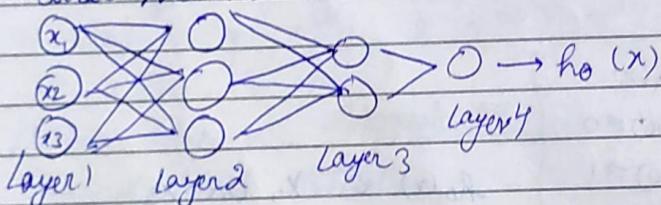


hidden layer

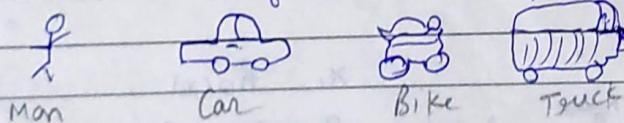
$x_1 \cdot x_2 \quad a_1 \quad a_2 \quad h_0(x)$

$(0, 0)$	0	0	1	1	AND gate
$(0, 1)$	0	1	0	0	XNOR gate
$(1, 0)$	1	0	0	0	using diff. gates.
$(1, 1)$	1	1	0	1	

⇒ Neural Network Intuition



⇒ Multiclass Classification = Multiple output units one Vs all



0	0	$h_0(x) \in \mathbb{R}^4$
0	0	$\rightarrow \text{man}$
0	0	$\rightarrow \text{car}$
0	0	$\rightarrow \text{bike}$
0	0	$\rightarrow \text{truck}$

Output will be
a vector of 4 numbers

Output for man will be $h_0(x) \in \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$; $h_0(x) \in \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$

$$h_0(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Man

car

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

$y^{(i)}$ one of $\{ \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \}$

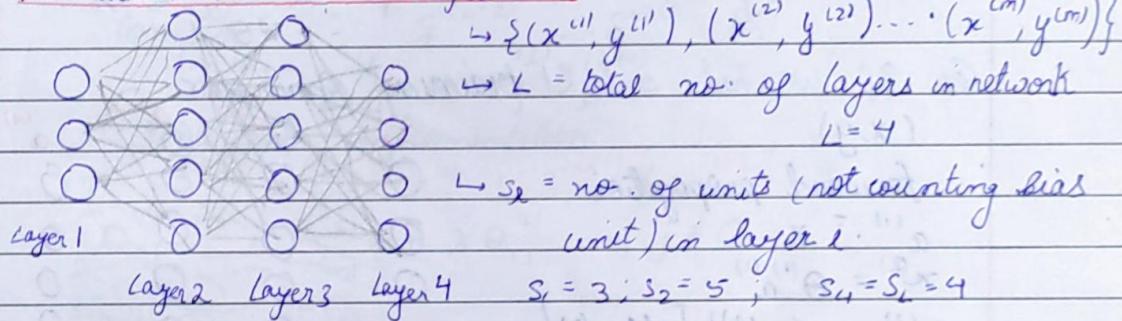
$x^{(i)}$ man car Bike Truck

$(x^{(i)}, y^{(i)})$ $\underbrace{h_0(x^{(i)}) \in y^{(i)}}_{\text{IR}}$

Previously
 $y \in \{1, 2, 3, 4\}$

Week - 5

* Neural Networks - Cost function:



Binary classification

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

1 output unit

$$h_0(x) \in \text{IR}$$

$$S_L = 1; K = 2$$

multi class - Classification (K classes)

$$y \in \text{IR}^K \text{ & } \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

man car Bike Truck

K output units

$$h_0(x) \in \text{IR}^K$$

$$S_2 = K \quad (K \geq 3)$$

\Rightarrow cost function:

logistic regression:

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_0(x^{(i)})) + (1-y^{(i)}) \log(1-h_0(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$$H_1 = \gamma_1 w_1 + \gamma_2 w_2 + b_1$$

$$\xi = \frac{1}{1+e^{-\xi}} \rightarrow \eta_1 = \xi$$

~~$$H_1 = \gamma_1 w_1 + \gamma_2 w_2 + b_1$$~~
~~$$H_2 = \gamma_1 w_1 + \gamma_2 w_2 + b_2$$~~

Page No. 40
Date _____

Neural Networks :-

$\text{Ro}(x) \in \mathbb{R}^k; (\text{Ro}(x))_i = i^{\text{th}} \text{ Output}$

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

Sum of $y_k^{(i)}$
 over $y_k^{(i)}$
 $y_k^{(i)}$'s

$\theta_{j|i}^{(l)} \in \mathbb{R}$

⇒ Backpropagation algorithm :-

so we need to compute $J(\theta) \triangleq \frac{\partial}{\partial \theta_j^{(l)}} J(\theta)$

Ex If we had only 1 training example :-

(x, y)

∴ forward propagation

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

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

$$a^{(2)} = g(z^{(2)}) \quad [\text{add } a_0^{(2)}]$$

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

$$a^{(3)} = g(z^{(3)}) \quad [\text{add } a_0^{(3)}]$$

$$z^{(4)} = \theta^{(3)} a^{(3)}$$

$$a^{(4)} = \text{Ro}(x) = g(z^{(4)})$$

$$\begin{array}{cccc}
 a^{(1)} & 0^{(2)} & G^{(3)} & a^{(4)} \\
 \downarrow & \downarrow & \downarrow & \downarrow \\
 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0
 \end{array}$$

⇒ Compute gradient using Backpropagation algorithm :-

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

For each output unit ($\text{layer } L=4$)

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad \text{or} \quad S^{(4)} = a^{(4)} - y$$

$\hookrightarrow (\text{Ro}(x))_j$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \cdot g'(x^{(3)})$$

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

No ($\delta^{(1)}$) as these are features so no errors in it

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \cdot s_i^{(l+1)} \quad [\text{ignoring } \lambda; \delta_0 = 0]$$

\Rightarrow Backpropagation algorithms :-

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

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

and to compute $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

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

Set $a^{(1)} = x^{(1)}$

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

(Using $y^{(1)}$, compute $s^{(L)} = a^{(L)} - y^{(L)}$)

Compute $s^{(L-1)}, s^{(L-2)}, \dots, s^{(2)}$

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

(vectorized implementation)

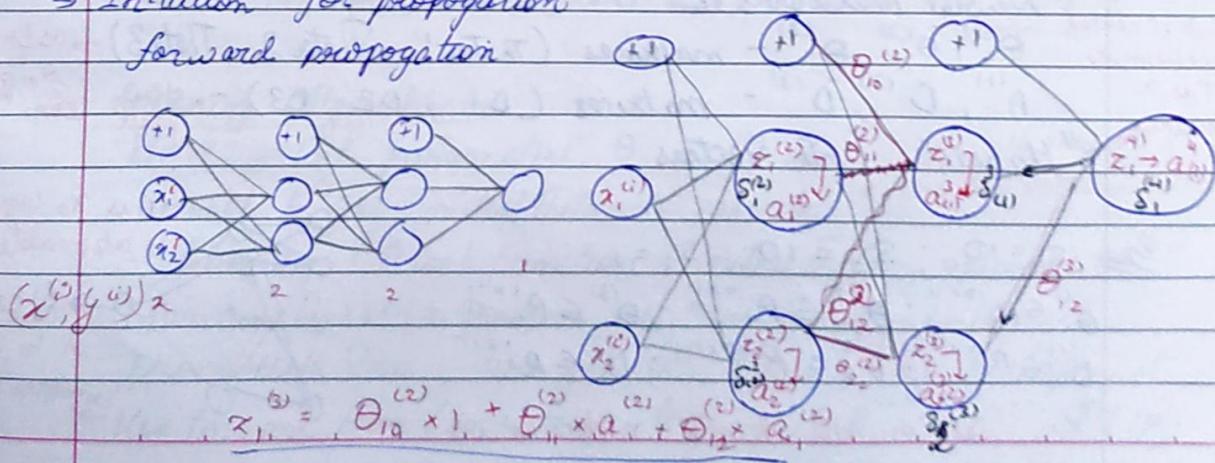
$\Delta^{(1)} := \Delta^{(1)} + s^{(2)} (a^{(1)})^T$

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

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

$\therefore \frac{\partial}{\partial \theta_{ij}^{(1)}} J(\theta) = D_{ij}^{(1)}$

\Rightarrow Intuition for propagation :-
forward propagation



Now for single example $x^{(i)}, y^{(i)}$, the case of 1 output unit and ignoring regularization ($\lambda=0$)

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

$$\text{This cost}(i) \approx (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

i.e how well is the network doing on example i?

for backward propagation

$$\delta_j^{(l)} = \text{"Error" of cost for } a_j^{(l)} \text{ (unit } j \text{ in layer } l)$$

Formally $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$ (for $j \geq 0$), where

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

like $\delta_1^{(4)} = y^{(i)} - Q^{(4)}$
product - value predicted

$$\delta_2^{(2)} = \Theta_{12}^{(2)} \cdot \delta_{(1)}^{(1)} + \Theta_{22}^{(2)} \cdot \delta_{(2)}^{(3)}$$

$$\delta_2^{(3)} = \Theta_{12}^{(3)} \cdot \delta_{(1)}^{(1)}$$

* Backpropagation implementation - Unrolling parameters:
Advanced Optimization:

function [$\text{Val}, \text{gradient}$] = costFunction(theta)

In logistic regression

$\hookrightarrow \mathbb{R}^{n+1}$

$\hookrightarrow \mathbb{R}^{n+1}$ (vectors)

theta was optTheta = fminunc(@(costFunction, initialTheta, options)

a vector

neural Networks ($L=4$):

but it neural networks it as a matrix $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices ($\Theta_1, \Theta_2, \Theta_3$) \rightarrow as $L=4$

$D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D_1, D_2, D_3) \rightarrow new grad term

"Unroll" into vectors

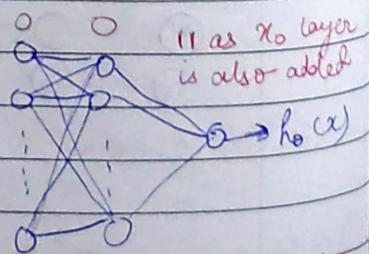
$$\text{Ex} \quad S_1 = 10; S_2 = 10; S_3 = 1$$

$$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}; \Theta^{(2)} \in \mathbb{R}^{10 \times 11}; \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

$$D^{(1)} \in \mathbb{R}^{10 \times 11}; D^{(2)} \in \mathbb{R}^{10 \times 11}; D^{(3)} \in \mathbb{R}^{1 \times 11}$$

let convert it into matrices in

Octave



$$\Theta^{(1)} \quad \Theta^{(2)} \quad \Theta^{(3)}$$

$$\text{Theta Vec} = [\Theta^{(1)}(:); \Theta^{(2)}(:); \Theta^{(3)}(:)];$$

$$D \text{Vec} = [D^{(1)}(:); D^{(2)}(:); D^{(3)}(:)];$$

to give values
3 reshaped
e...
 $\Theta^{(1)} = \text{reshape}(\text{ThetaVec}(1:110), 10, 11);$
 $\Theta^{(2)} = \text{reshape}(\text{ThetaVec}(111:220), 10, 11);$
 $\Theta^{(3)} = \text{reshape}(\text{ThetaVec}(221:231), 1, 11);$

Let's code in octave.

$$\Theta^{(1)} = \text{ones}(10, 11) \rightarrow \Theta^{(1)} = \begin{matrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{matrix} \quad \left\{ \begin{matrix} 10 \\ 11 \end{matrix} \right\}$$

$$\Theta^{(2)} = 2 * \text{ones}(10, 11) \rightarrow \Theta^{(2)} = \begin{matrix} 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \end{matrix} \quad \left\{ \begin{matrix} 10 \\ 11 \end{matrix} \right\}$$

$$\Theta^{(3)} = 3 * \text{ones}(1, 11) \rightarrow \Theta^{(3)} = \begin{matrix} 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 \end{matrix} \quad \left\{ \begin{matrix} 1 \\ 11 \end{matrix} \right\}$$

* To min & make thetaVec

$$\text{thetaVec} = [\Theta^{(1)}(:); \Theta^{(2)}(:); \Theta^{(3)}(:)] \quad \text{thetaVec} = \begin{matrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 \end{matrix} \quad \left\{ \begin{matrix} 110 \\ 11 \\ 11 \end{matrix} \right\}$$

To extract or reshape thetaVec to get our theta

$$\Theta^{(1)} = \text{reshape}(\text{thetaVec}(1:110), 10, 11) \rightarrow \Theta^{(1)} = \begin{matrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{matrix} \quad \left\{ \begin{matrix} 10 \\ 11 \end{matrix} \right\}$$

$$\Theta^{(2)} = \text{reshape}(\text{thetaVec}(111:220), 10, 11) \rightarrow \Theta^{(2)} = \begin{matrix} 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \end{matrix} \quad \left\{ \begin{matrix} 10 \\ 11 \end{matrix} \right\}$$

$$\Theta^{(3)} = \text{reshape}(\text{thetaVec}(221:231), 1, 11) \rightarrow \Theta^{(3)} = \begin{matrix} 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 & 3 \end{matrix} \quad \left\{ \begin{matrix} 11 \\ 1 \end{matrix} \right\}$$

2) Learning Algorithm

Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ we will pass $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ into long matrix initialTheta

Unroll to get initialTheta to pass to fminunc (Qcost function, initialTheta, options)

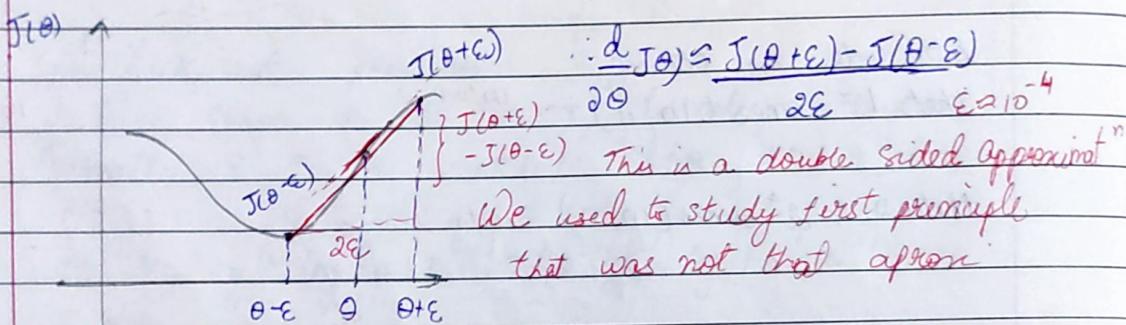
function [J, grad] = costFunction(thetaVec)

From thetaVec we get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ by reshaping it

Use forward prop/back prop to compute $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ & J(theta)

Unroll $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ to get gradientVec → then we will find it

⇒ Gradient checking : what was the need of gradient checking because when we use forward/back propagation like for 8 calculation there are chances of bug or errors. So we use grad checking.
Numerical estimation of gradients



Implementation :

$$\text{gradApprox} = \frac{(J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON}))}{(2 * \text{EPSILON})}$$

⇒ Parameter vector θ :

$\theta \in \mathbb{R}^n$ (E.g. θ is "unrolled" version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

$$\theta = \theta_1, \theta_2, \theta_3, \dots, \theta_n$$

$$\frac{\partial J(\theta)}{\partial \theta_1} \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial J(\theta)}{\partial \theta_2} \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$$

⋮

$$\frac{\partial J(\theta)}{\partial \theta_n} \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_{n-1}, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_{n-1}, \theta_n - \epsilon)}{2\epsilon}$$

Let's implement it in octave :

for $i = 1:n$,
 thetaPlus = theta;
 $\{\theta_1, \theta_2, \dots, \theta_n\}$
 $\{\theta_1 + \epsilon, \theta_2, \dots, \theta_n\}$
 $\{\theta_1 - \epsilon, \theta_2, \dots, \theta_n\}$
 thetaMinus = theta;
 $\{\theta_1 - \epsilon, \theta_2, \dots, \theta_n\}$
 $\text{gradApprox}(i) = (\mathcal{J}(\theta_{\text{plus}}) - \mathcal{J}(\theta_{\text{minus}})) / (2 * \epsilon)$;
 end;
 $\mathcal{J} \approx \mathcal{J}(\theta)$
 Now check that $\text{gradApprox} \approx \Delta \text{Vec}$ this we got from back propagation

⇒ Implementation Note:

- i) Implement backprop to compute ΔVec (unrolled $\Delta^{(1)}, \Delta^{(2)}, \Delta^{(3)}$)
- ii) Implement numerical gradient check to compute gradApprox
- iii) Make sure they give similar values
- iv) Turn off gradient checking. Using backprop code for learning.

⇒ Important:

- i) Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of costFunction(...)) your code will be very slow.

⇒ Random initialization:

Initial value of θ

For gradient descent and advanced optimization method, need initial value for θ .

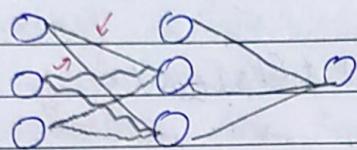
`optTheta = fminunc (CostFunction, initialTheta, options)`

So can we consider gradient descent X Don't do that
Set `initialTheta = zeros(n, 1)`? X in neural network

this works fine

for regressions but
not here

If we initialize θ as zero



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

$a_1^{(2)}$, $a_2^{(2)}$ will be same

$$\therefore a_1^{(2)} = a_2^{(2)}$$

Also $s_1^{(2)} = s_2^{(2)}$

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

$$\therefore \theta_{01}^{(1)} = \theta_{02}^{(1)}$$

Non zero maybe
But same

Some constant C

After each update, parameters corresponding to inputs going into each of two hidden units are identical.

$$a_1^{(2)} = a_2^{(2)} \text{ i.e. all hidden units are same}$$

ie. initialized, to - can

We use random initialization : Symmetry breaking
Initialize each $\theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$

$$(i.e. -\epsilon \leq \theta_{ij}^{(l)} \leq \epsilon)$$

Eg:- random 10×11 matrix (10×0.8)

$[-\epsilon, \epsilon]$

$$\Theta_1 = \text{rand}(10, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

$$\Theta_2 = \text{rand}(1, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

Putting it together :-

Pick a neural architecture (connectivity pattern between neurons)

3 5 4
Hidden layer

3 5 5 4
Hidden layers

3 5 5 5 4
Hidden layer

No. of input units = Dimensions of features $x^{(0)}$

No. of output units = Number of classes

$$y \in \{1, 2, 3, \dots, 10\}, \quad y = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix}$$

Reasonable default : 1 hidden layer, or if > 1 hidden layer, have same no. of hidden units in every layer (usually the more the better).

- ⇒ Let's train a neural network -
- i) Randomly initialize weights
- ii) Implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
- iii) Implement code to compute cost function $J(\theta)$
- iv) Implement back prop to compute partial derivatives $\frac{\partial J(\theta)}{\partial \theta_j}$.

for $i = 1:m$, $\{ (x^{(i)}, y^{(i)}) | (x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)}) \}$

Perform forward propagation and back propagation using example $(x^{(i)}, y^{(i)})$.

(Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l=2, 3, \dots, L$).

$$\{ \Delta^{(l)} := A^{(l)} + S^{(l+1)} (a^{(l)})^T \dots \}$$

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

v) Use gradient checking to compare $\frac{\partial}{\partial \theta_j} J(\theta)$ computed using backprop vs. using numerical estimate of gradient of $J(\theta)$. Then disable gradient checking code.

vi) Use gradient descent or advance optimization method with backpropagation to try to minimize $J(\theta)$ as a "function" of parameters θ . $\rightarrow \frac{\partial}{\partial \theta_j} J(\theta)$ Non-convex \rightarrow Not a big problem $J(\theta)$

* Neural networks on graph

* Autonomous Driving

Week - 6* Advice for applying Machine Learning

Debugging a learning algorithm :

Suppose you have implemented regularized linear regression to predict housing prices.

$$T(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)}) - y^{(i)} \right]^2 + \lambda \sum_{j=1}^m \theta_j^2$$

However when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- ↪ Get more training examples
- ↪ Try smaller sets of features
- ↪ Try getting additional features
- ↪ Try adding polynomial features ($x_1^2, x_2^2, x_1 x_2$ etc.)
- ↪ Try decreasing λ
- ↪ Try increasing λ

⇒ Machine learning Diagnostics :-

Diagnostics : A test that you can run to gain insight what is / isn't working with a learning algorithm and gain guidance as to how best to improve its performance.

Diagnostics can take time to implement, but doing so can be a very good use of your time.

⇒ Evaluation of hypothesis :-

We know about overfitting of data that is one issue as the overfit data completely plots the graph but fails to predict test data i.e. not in training set.

⇒ Evaluating your hypothesis :-
Dataset

Size	Price		
2104	400		$(x^{(1)}, y^{(1)})$
1600	330	Training set	$(x^{(2)}, y^{(2)})$
2400	369		
1416	232		
3000	540		
1985	300		
1534	315		$(x^{(m)}, y^{(m)})$
1427	199	Test set	$(x_{\text{test}}, y_{\text{test}})$
309	1380	set	$(x_{\text{test}}, y_{\text{test}})$
1490	243		$(x_{\text{test}}, y_{\text{test}})$

⇒ Training / testing procedure for linear regression :-

↳ Learn parameters θ from training data (minimize $J_{\text{train}}(\theta)$)

↳ Compute test set error - $J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (\hat{y}_i - y_i)^2$ $\approx 70\%$.

→ for logistic regress?

$$J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_i \ln \hat{y}_i + (1 - y_i) \ln (1 - \hat{y}_i)$$

Misclassification error (0/1 misclassification error)

$$\text{err}(\hat{y}_i, y_i) = \begin{cases} 1 & \text{if } \hat{y}_i > 0.5, \\ & \quad y=0 \\ 0 & \text{otherwise} \\ & \quad y=1 \end{cases}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(\hat{y}_i, y_i)$$

we use $J_{\text{val}}(\theta)$
here
and then put it for
 J_{test}

⇒ Model Selection :- $d = \text{degree of polynomial}$ $\rightarrow J_{\text{val}}(\theta)$

$$\left\{ \begin{array}{l} d=1 \\ \hat{y}(x) = \theta_0 + \theta_1 x \end{array} \right. \rightarrow \theta^{(1)} \text{ parameters i get by fitting this model}$$

$$\left\{ \begin{array}{l} d=2 \\ \hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \end{array} \right. \rightarrow \theta^{(2)} \rightarrow J_{\text{test}}(\theta^{(2)})$$

if any $J_{\text{val}}(\theta^{(d)})$ is min we optimistically estimate of generalization error

$$\left\{ \begin{array}{l} d=10 \\ \hat{y}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \end{array} \right. \rightarrow \theta^{(10)} \rightarrow J_{\text{test}}(\theta^{(10)})$$

⇒ Evaluating your hypothesis :

	Size	Price	
	2104	400	Training ($x^{(i)}, y^{(i)}$)
	1600	330	set
60%	2400	369	
	1416	232	
	3000	540	
	1985	300	($x^{(m)}, y^{(m)}$)
	1534	315	Cross Validation : ($x_{cv}^{(i)}, y_{cv}^{(i)}$)
20%	1427	199	Set (CV) ($x_{cv}^{(mcv)}, y_{cv}^{(mcv)}$)
	1380	212	Test
20%	1494	243	Set ($x_{test}^{(mcv)}, y_{test}^{(mcv)}$)

$m_{cv} = \text{no. of CV}$

Example : ($x_{cv}^{(i)}, y_{cv}^{(i)}$)

⇒ Train / validate / test error :

Training Error :

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

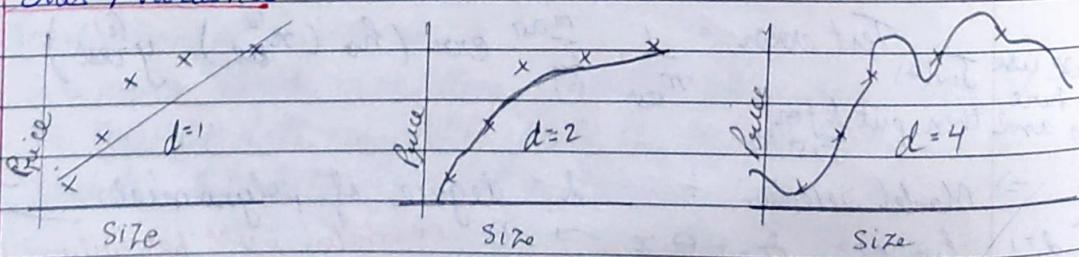
Cross Validation Error :

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test Error :

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

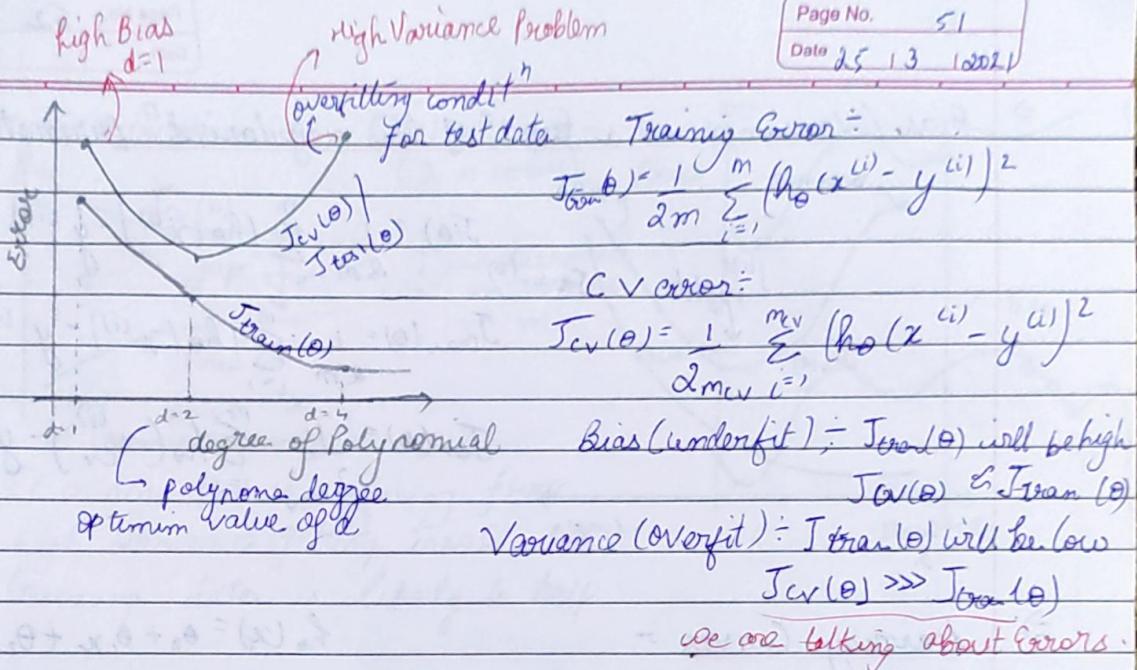
⇒ Bias / Variance



$\theta_0 + \theta_1 x$
High Bias
(Underfit)

$\theta_0 + \theta_1 x + \theta_2 x^2$
"Just Right"

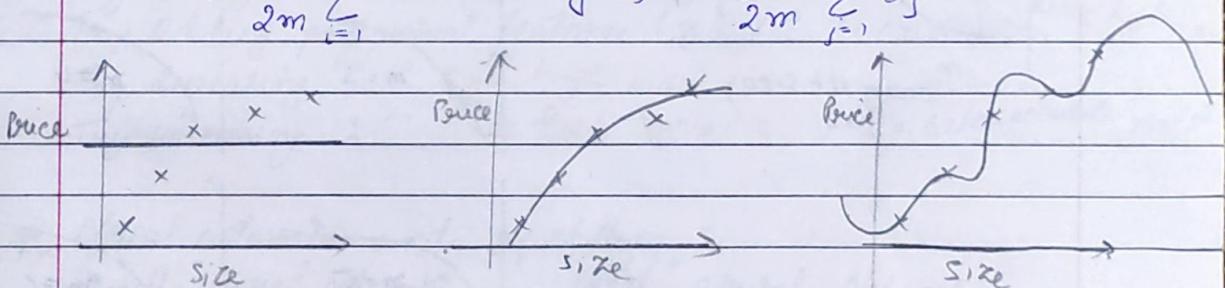
$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
High Variance
(Overfit)



3) Impact of Regularization on Variance & Bias:

Model: $\hat{\theta}_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\hat{\theta}_0(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$



large $\lambda \approx 1000$

High Bias (underfit)

$\theta_1 \approx 0, \theta_2 \approx 0 \dots$

$\hat{\theta}_0(x) \approx \theta_0$

Intermediate λ

"Just right"

Small $\lambda \approx 0$

High Variance (overfit)

It will follow 4 degree polynomial

How to choose the regularization parameter λ ?
we can try diff. values of λ :

say $\lambda = 0$

$\rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$

$\lambda = 0.01 \left\{ \begin{array}{l} \rightarrow \min_{\theta} (J(\theta)) \rightarrow \theta^{(2)} \rightarrow J_{\text{cv}}(\theta^{(2)}) \rightarrow J_{\text{cv}}(\theta^{(5)}) \\ \rightarrow \theta^{(3)} \end{array} \right.$

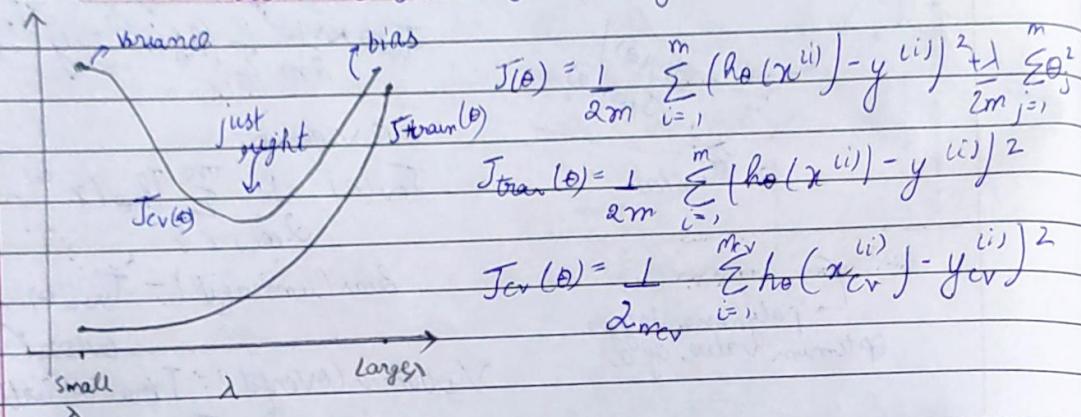
$\lambda = 0.02 \left\{ \begin{array}{l} \rightarrow \min_{\theta} (J(\theta)) \rightarrow \theta^{(4)} \rightarrow J_{\text{cv}}(\theta^{(4)}) \\ \rightarrow \theta^{(5)} \end{array} \right.$

$\lambda = 10 \left\{ \begin{array}{l} \rightarrow \min_{\theta} (J(\theta)) \rightarrow \theta^{(6)} \rightarrow J_{\text{cv}}(\theta^{(6)}) \\ \rightarrow \theta^{(7)} \end{array} \right.$

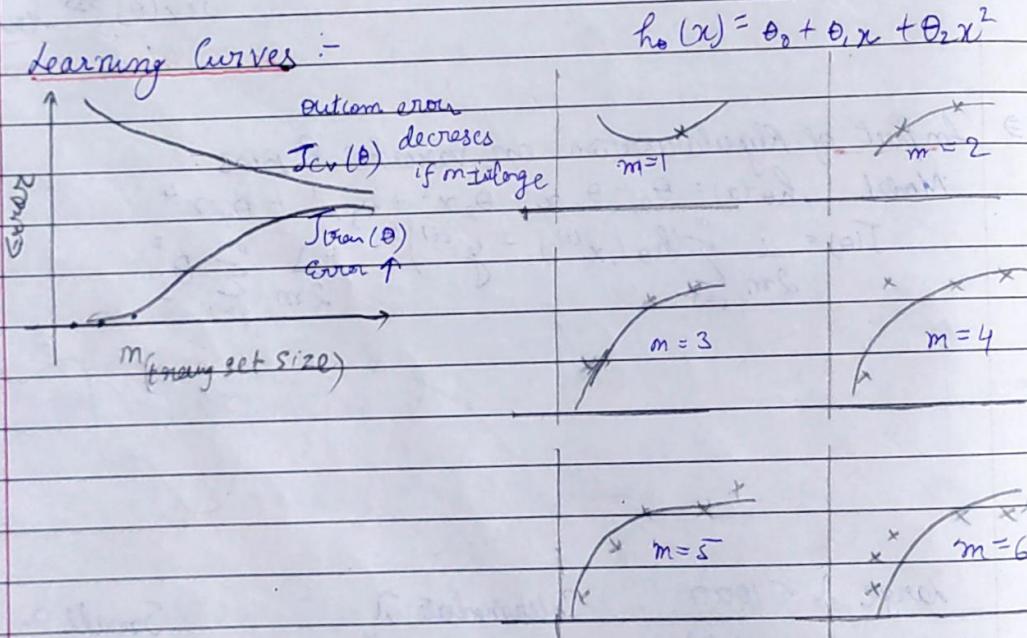
$\rightarrow \theta^{(12)} \rightarrow J_{\text{cv}}(\theta^{(12)})$

Pick. (say) $\theta^{(5)}$. Test Error = $J_{\text{test}}(\theta^{(5)})$

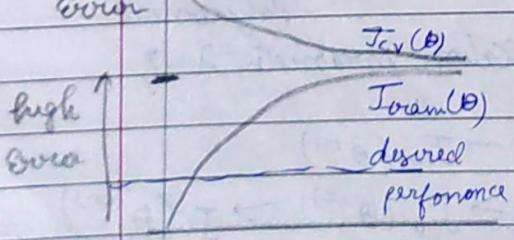
⇒ Bias / Variance as a functⁿ of the regularizatⁿ parameter λ :-



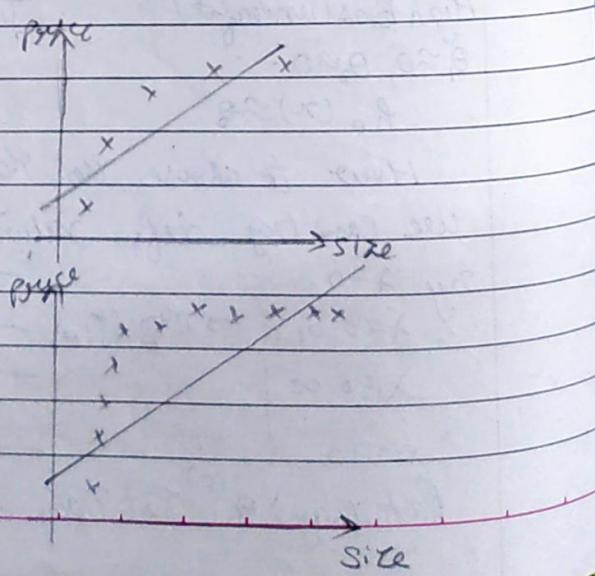
⇒ Learning Curves :-



⇒ High Bias $= h_\theta(x) = \theta_0 + \theta_1 x$



If a algorithm is suffering from high bias ; getting more training data will not help much by itself.

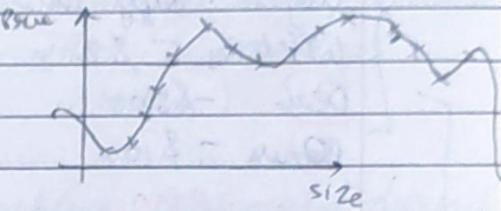
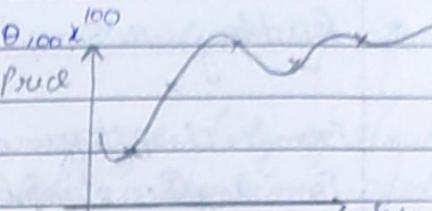


⇒ High Variance : $R_{\theta}(x) = \theta_0 + \theta_1 x - \cdot \cdot \cdot \theta_{100} x^{100}$
 $(\lambda \text{ is small})$ $J_{CV}(\theta)$

Desired performance
--- gap --- ↑ gap decreases
on increase in m

$J_{train}(\theta)$
 m (training set size)

If a algorithm is suffering from high variance getting more training data is likely to help.



⇒ Debugging a learning algorithm :

Get more training Examples → fixes high variance

Try smaller sets of features → fixes high variance

Try getting additional features → fixes high bias

Try adding polynomial features ($x_1^2, x_2^2, x_1 x_2$ etc) → fixes high bias

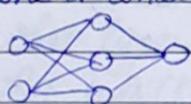
Try decreasing λ → fixes high bias

Try increasing λ → fixes high variance for a better accurate model

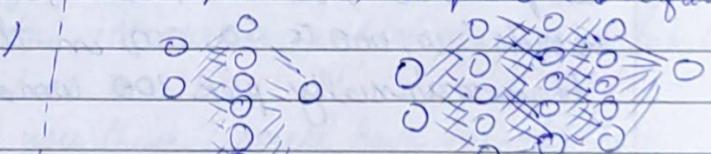
⇒ Neural networks and overfitting :

"Small" neural network ; "large" neural network

(fewer parameters ; more parameters ; more prone to overfitting)
prone to underfitting



computational cheaper



computational more expensive
use regularization (λ) to address overfitting

* Building a spam classifier

From: cheapstuffs.com

To: any@cs.stanford.edu

Subject: ~~Buy now!~~

Watches - \$2

Deal - \$500

Now - \$100

spam

From: Alfred Ng

To: any@cs.stanford.edu

Subject: Meet

Hey Andrew,

Let's meet tomorrow.

Alf

Non-Spam

⇒ Build a spam classifier :-

Supervised learning x = features of email

y = spam (1) or not spam (0)

Features x : choose 100 words indicative of spam/not spam

e.g. deal, buy, discount, andrew, now ...

Now $x = \{x_j\}_{j=1}^{100}$ andrew $x_j = 1$ if word j appears

$x \in \mathbb{R}^{100}$ buy $x_j = 0$ in email

deal $x_j = 0$ otherwise

discount

now

Alphabetical order

In practice, take most frequently occurring n words (10,000 to 50,000) in training set, rather than manually pick 100 words

⇒ Building a spam classifier

How to spend your time to make it have low error?

→ collect lots of data

→ e.g. "honeypot" project

→ Develop sophisticated features based on email routing information (from email header).

- ↳ Develop sophisticated features based on email routing information (from email header).
 - ↳ Develop sophisticated features for message body e.g. should "discount" & "discounts" be treated as the same word? How about "deal" and "Dealer"? Features about punctuation.
 - ↳ Develop sophisticated algorithm to detect misspellings (e.g. ~~mortgage~~, med. lincne, w4tches)
- ⇒ Error Analysis :- How to start a problem ; Recommended Approach
- ↳ Start with a simple algorithm that you can implement quickly. Implement it and test it on your cross validation data.
 - ↳ Plot learning curves to decide if more data ; more features etc are likely to help.
 - ↳ Error Analysis - Manually examine the examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of example it is making errors on.
- ⇒ Error Analysis :-
- $n_{cv} = 500$ examples in cross validation set.
- Algorithm misclassifies 100 emails.
- Manually examine the 100 errors, and categorize them based on -
- i) What type of email it is → pharma, replica, steal passwords ...
 - ii) What cues/features you think would have helped the algorithm classify them correctly
- | | |
|---------------------|---------------------------------------|
| Pharma = 12 | ↳ Deliberate misspellings = 5 |
| Replica / fake = 4 | (mortgage, med. lincne, etc) |
| Steal password = 53 | ↳ Unusual email routing = 16 |
| Other = 31 | ↳ Unusual (spamming) punctual "(say)" |

⇒ The importance of numerical evaluation :-

↳ Should discount / discounts / discounted / discounting be treated as the same word?

↳ Can use "stemming" software (e.g. "Porter stemmer")
used in NLP (natural language processing) software "stemming"

can be sometimes difficult to differentiate

universe / universality: same starting initials

error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try it and see if it works.

Need Numerical evaluation (e.g. cross validation score) of algorithm's performance with and without stemming

without stemming

5% error

with stemming

3% error

Distinguish upper v/s lower case (Morn / mom) 3.2% error

⇒ Error metrics of skewed classes :-

Cancer classification example :-

Train logistic regression model $h_0(x)$. If $y=1$ if cancer
Find that you got 1% error on test set $y=0$ otherwise
(99% correct diagnosis)

Only 0.5% of patients have cancer. Now 1% error is huge
function $y = \text{predictCancer}(x)$ ↗ skewed classes
 $y=0$ ↗ ignore x ↗ This will give only
return. 0.5% error.

⇒ Precision / Recall :-

$y=1$ in presence of rare class that we want to detect
Actual Class Precis :- [of all patients where we predicted $y=1$ what fraction actually has cancer]

Predicted class	1	0
1	True Positive	False Positive
0	False Negative	True Negative

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

⇒ Recall (of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

True positives = True positives

Factual positives = True pos + False negatives

If our model predicts $y=0$ always \therefore recall = 0

Model with high precⁿ & recall is a good model.

⇒ Trading off precision and recall:

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

Predict 1 if $h_{\theta}(x) \geq 0.5$ 0.7 → High precⁿ, lower recall

Predict 0 if $h_{\theta}(x) < 0.5$ 0.7 ↙

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

ii) Suppose we want to avoid missing too many cases of cancer (avoid false negatives)

∴ we can keep the threshold low

1 $h_{\theta}(x) \geq 0.5$ 0.3 → higher recall, lower precision

0 $h_{\theta}(x) < 0.5$ 0.3 ↗

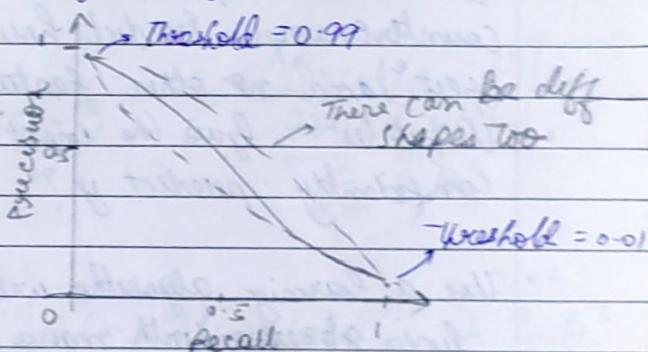
More generally predict 1 if $h_{\theta}(x) \geq$ threshold

Precⁿ = true positives

No. of predicted positives

Recall = True positives

No. of actual positives



⇒ F-score (F score) =

How to compare precision / recall numbers?

	Precision (P)	Recall (R)	Average	F-score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.51	0.0392

$$F\text{-score} = \frac{2PR}{P+R} \quad \text{if } P=0 \text{ or } R=0, \rightarrow F\text{-score} \geq 0 \\ P=1 \text{ or } R=1 \rightarrow F\text{-score} = 1$$

\Rightarrow Data for machine learning :-

Design a high accuracy learning system

Eg:- Classify between confusable words.
 {to, two, too } , {then, than}

For breakfast I ate ~~two~~ eggs.

Algorithms :-

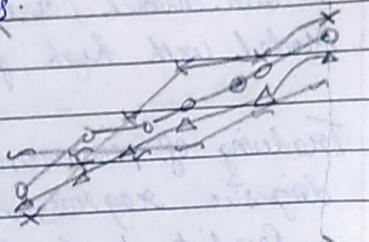
\hookrightarrow Perception [Logistic Regression]

\hookrightarrow Widnow

\hookrightarrow Memory-based

\hookrightarrow Naïve Bayes

It's not who has the best algorithm
 that wins. It's who has the
 most data wins!



Training set size 1000
 in millions

\Rightarrow Large data rationale :-

\hookrightarrow Assume feature $x \in \mathbb{R}^{n+1}$ has sufficient information to predict y accurately.

\hookrightarrow For breakfast I ate ~~two~~ eggs. ~~Ex~~

Counterexample:- Predict housing price from only size (feet²) and no other features.

\hookrightarrow Useful test:- Give the input x , can a human expert confidently predict y ?

\hookrightarrow Use a learning algorithm with many parameters (e.g. logistic regression with many features, neural network with many hidden units). ~~also this algorithm~~

(many error) $J_{\text{train}}(\theta)$ will be small

use a very large training set (unlikely to overfit) ~~and variance~~

$J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

$J_{\text{test}}(\theta)$ will be small.

Week - 7

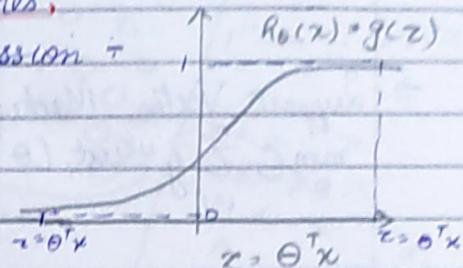
Support Vector Machines.

\Rightarrow Alternate view of logistic regression :-

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

If $y=1$; we want $h_{\theta}(x) \approx 1$
 $\Leftrightarrow \theta^T x > 0$

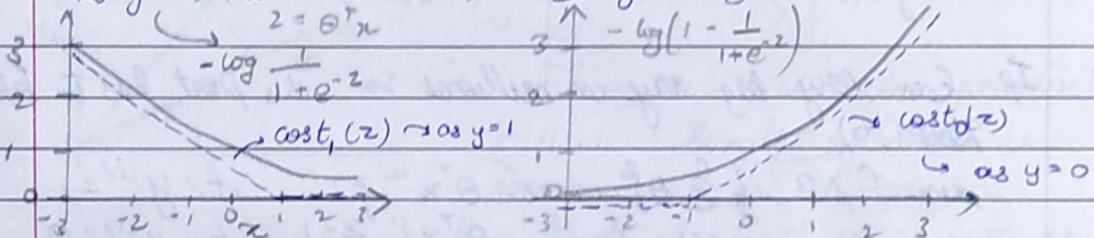
If $y=0$; we want $h_{\theta}(x) \approx 0$
 $\Leftrightarrow \theta^T x \ll 0$



$$\Rightarrow \text{Cost of example} : -(y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x)))$$

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

If $y=1$; (want $\theta^T x \gg 0$) If $y=0$; (want $\theta^T x \ll 0$)



Logistic Regression:

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

Support Vector Machine:

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_t(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_o(\theta^T x^{(i)}) + \frac{\lambda}{2m} \sum_{j=0}^n \theta_j^2$$

constant can be removed

$$\text{Ex } \min_u K[(u-s)^2 + 1] \rightarrow u=s \text{ : no use of K.}$$

$$A + \lambda B \rightarrow CA + B \quad ? \text{ Different Methods of parameterizing }$$

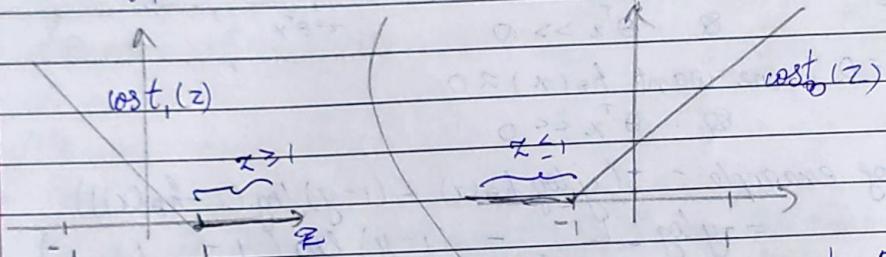
SVM :- $\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_t(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_o(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$

Hypothesis :- C A B

$$h_0(x) \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{Otherwise} \end{cases}$$

\Rightarrow Support Vector Machine:

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



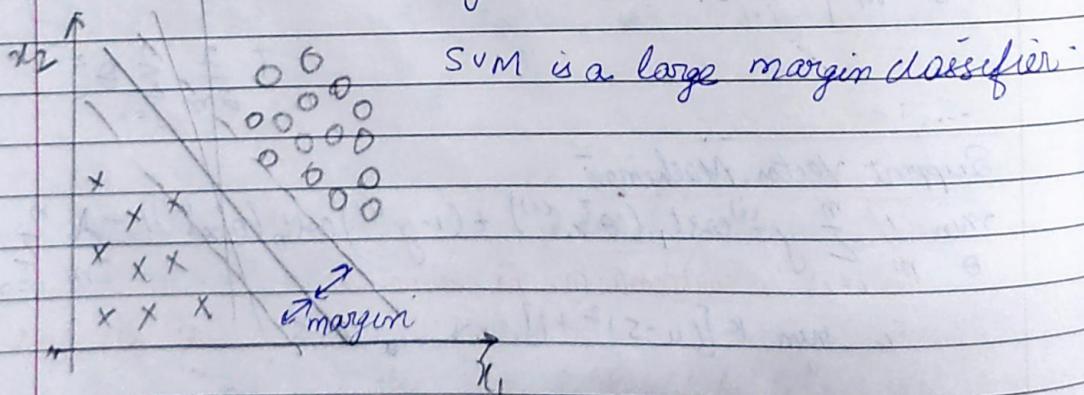
If $y=1$, $\theta^T x \geq 1$
(not just ≥ 0)

If $y=0$; we want $\theta^T x \leq -1$
(not just < 0)

If C is very big say in millions so this part has to be zero(0)

$$\min C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad \text{st: } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

\Rightarrow SVM Decision Boundary = Linear Separable Case:



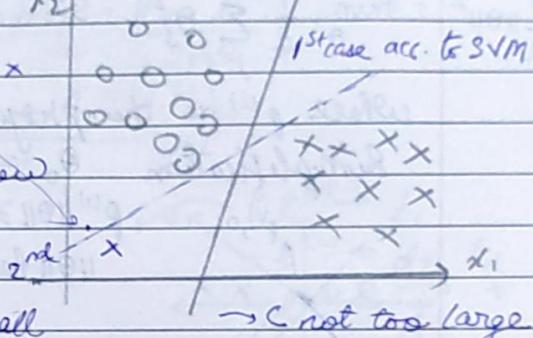
\Rightarrow Large margin classifier in presence of outliers :-

Case 1st when we have clustered $\circ \otimes \times$

(Case 2) nd when there is one (\times) so our new SVM will be 2nd line if C is very large.

$C \rightarrow \sqrt{2}$ i.e. λ is very small

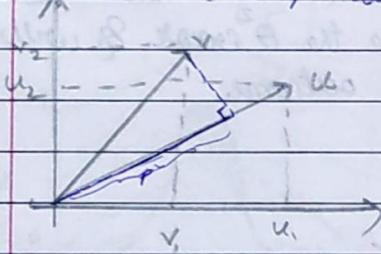
If C is not large we can ignore these small outliers.



$\rightarrow C$ not too large

\Rightarrow Maths behind SVM (optional but not for me lol) :-

Vector inner product :-



$$u = [u_1, u_2] ; v = [v_1, v_2] ; u^T = [u_1, u_2]$$

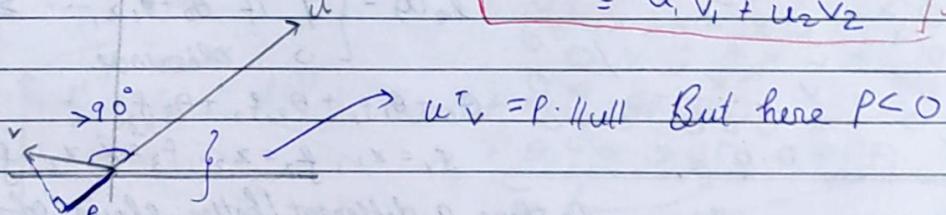
$u^T v$ is the vector inner product

$$\|u\| = \text{length of vector } u \\ \text{norme} = \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

Signed $\leftarrow P$ is length of projectⁿ of v onto w
it can be (+) or (-)

$$u^T v = P \cdot \|u\| = v^T u \quad \begin{cases} \text{Vector Inner} \\ \text{Product} \end{cases}$$

$$= u_1 v_1 + u_2 v_2$$



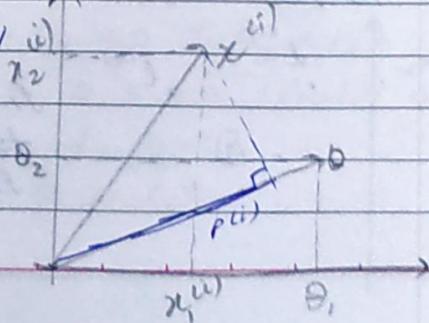
\Rightarrow SVM decision boundary :- $\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2)$
st. $\theta^T x^{(i)} \geq 1$ if $y^{(i)} = 1$ for simplification $\theta_0 = 0 \otimes n=2$

$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

$$\therefore \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \|\theta\| \quad \begin{matrix} \text{if} \\ \theta^T v \end{matrix}$$

$$\theta^T v = ? = P \cdot \|\theta\|$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$



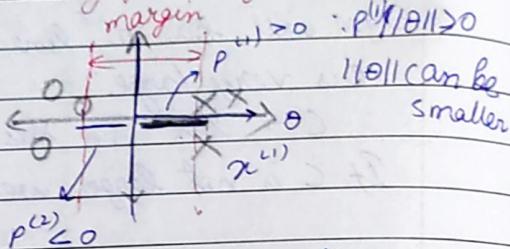
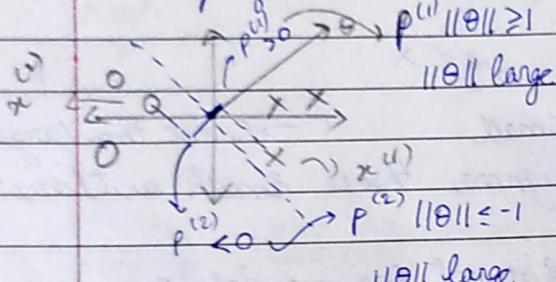
⇒ SVM Decision Boundary :-

$$\frac{1}{2} \|\theta\|^2 = \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad \text{S.T. } p^{(i)} \cdot \|\theta\| \geq 1 \quad \text{if } y^{(i)} = 1$$

$$p^{(i)} \cdot \|\theta\| \leq -1 \quad \text{if } y^{(i)} = 0$$

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ

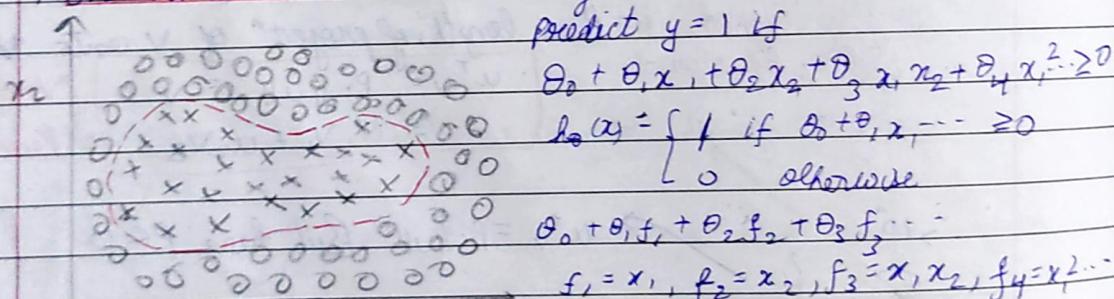
Simplification $\theta_0 = 0$:- decision boundary pass through origin.



This is the hypothesis of SVM
due to which θ will be small
which will reduce the θ^2 error & will
give best possible outcome.

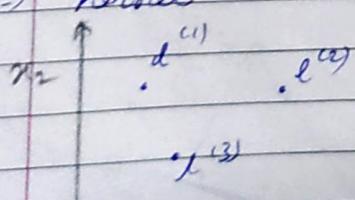
⇒ Kernels :-

Non-linear Decision Boundary :-



Is there a different/better choice of the feature f_1, f_2, f_3, \dots ?

⇒ Kernel :-



Given x , compute new features depending on proximity to landmarks.

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

$$f_2 = \dots$$

$$f_3 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

\Rightarrow Kernel

$$\text{Gaussian kernel } K(x, l^{(i)})$$

\Rightarrow Kernel & Similarity

$$f_1 = \text{Similarity } (x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$$

If $x \approx l^{(1)}$:

$$f_1 \approx \exp\left(-\frac{\sigma^2}{2\sigma^2}\right) \approx 1$$

$$l^{(1)} \rightarrow f_1$$

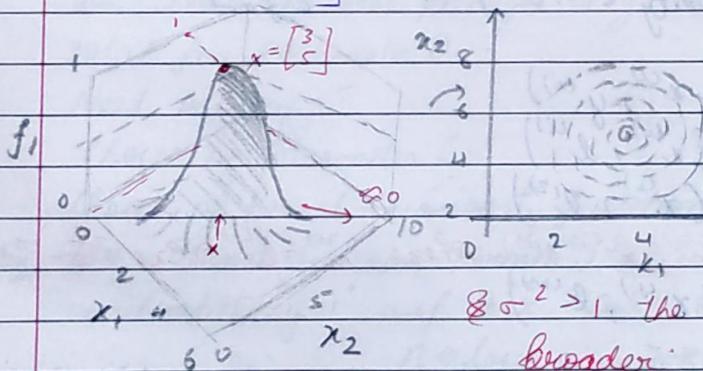
$$l^{(2)} \rightarrow f_2$$

$$l^{(3)} \rightarrow f_3$$

If x is far from $l^{(1)}$:

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

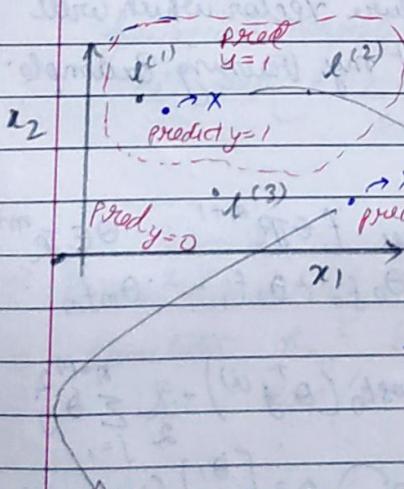
Ex $l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$; $f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) \approx \sigma^{-2} = 1$



If $\sigma^2 = 0$ the width of the bump will be narrower.

& contour plots will be close to each other.

If $\sigma^2 > 1$, the bell becomes broader and broader.



Predict "1" when:

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

$$x \rightarrow$$

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

As x is close to $l^{(1)}$

$$\therefore f_1 \approx 1; f_2 \approx 0 \& f_3 \approx 0$$

$$\theta_0 + \theta_1 \times 1 + \theta_2 \times 0 + \theta_3 \times 0 \approx$$

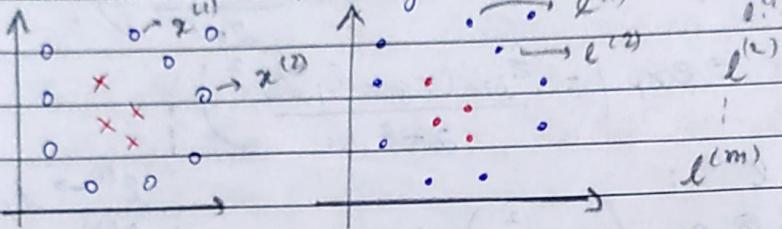
$$= -0.5 + 1 = 0.5 \geq 0 \quad \text{Predict } y = 1$$

As x is far from all $\therefore f_1, f_2, f_3 \approx 0$

$$\therefore \theta_0 + \theta_1 f_1 + \dots = \theta_0 = -0.5 < 0 \quad \text{predict } y = 0$$

\Rightarrow Kernels 2:

Now how can we get those $l^{(1)}, l^{(2)}, l^{(3)}, \dots, l^{(m)}$?



SVM with Kernels:

Given $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$, ..., $(x^{(m)}, y^{(m)})$

choose $l^{(1)} = x^{(1)}$; $l^{(2)} = x^{(2)}$; ..., $l^{(m)} = x^{(m)}$

Given example x :

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

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

...

$$f = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix}$$

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

$$x^{(i)} \rightarrow \begin{cases} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)}) \end{cases}$$

$$\therefore f^{(i)} = \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \quad f_0^{(i)} = 1 \quad \left\{ \begin{array}{l} \text{New feature vector which will} \\ \text{represent my training Example.} \end{array} \right.$$

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

Predict "y=1" if $\theta \cdot f \geq 0 \rightarrow \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$

Training:

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

$$\text{We can also write } \sum_j \theta_j^2 = \theta \cdot \theta \text{ on } \theta^T M \theta \quad \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \text{ (ignore } \theta_0\text{)} \quad \rightarrow \|\theta\|^2$$

\Rightarrow SVM parameters: weight

$c (= \lambda_2)$ Large c : lower bias, high variance (small λ)

~~underfit~~ Small c : Higher Bias, low Variance (large λ)

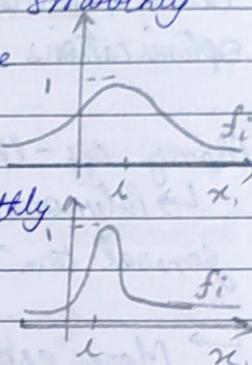
σ^2 large σ^2 : Features fit vary more smoothly

Higher bias, lower variance

$$\exp\left(-\frac{\|x - \mathbf{e}^{(0)}\|^2}{2\sigma^2}\right)$$

Small σ^2 : Features fit vary less smoothly.

lower bias, higher variance



\Rightarrow Using an SVM:

Use SVM softwares package (e.g liblinear, libsvm, ...) to solve for parameter θ .

Need to Specify :

choice of Parameter C

Choice of kernel (similarity function)

Ex No Kernel ("linear Kernel")

Predict " $y=1$ " if $\theta^T x \geq 0$ $\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$
 $n \rightarrow \text{large}$; $m \rightarrow \text{small}$ $x \in \mathbb{R}^{n+1}$

Gaussian Kernel :

$$f_i = \frac{\exp(-\|x - x^{(i)}\|^2)}{2\pi^2} \text{ where } x^{(i)} = x^{(i)}. \quad x \in \mathbb{R}^n \text{ and } m \text{ large}$$

Need to choose -²
use gaussian Kernel

A hand-drawn diagram of a cell with a nucleus, labeled "m > large".

⇒ Kernel (similarity) functions :

function $f_{\text{kernel}} = \text{Kernel}(x_1, x_2)$

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right) \quad x \in \mathbb{R}^2$$

return

f_m

Note: Do perform feature Scaling before using the Gaussian kernel

$$\|x - l\|^2 \quad \therefore \quad v = x - l \quad \& \quad \|v\| = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2} = (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2$$

$\sum_{i=1}^n x_i = \sum_{i=1}^n v_1 + \dots + v_n = v_1 + v_2 + \dots + v_n = (x_1, y_1) + (x_2, y_2) + \dots + (x_n, y_n)$

 \hookrightarrow 1-5 bedroom

$(x_1 - \bar{x})^2$ huge \therefore use feature scaling

⇒ Other choices of Kernel :

Note : Not all similarity functions similarity (x, i) make valid kernels. (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages' optimizations run correctly, and do not diverge)

⇒ Many off-the-shelf kernels available

↳ Polynomial Kernel : $K(x, l) = (x^T l)^2, (x^T l)^3, (x^T l)^4, \dots$
 General term = $(x^T l + \text{const})^{\text{degree}}$

↳ More esoteric : String Kernel, chi-square Kernel, Histogram, intersection Kernel.

similarity (x, l) if x is string we might use string kernel

⇒ Multi Class Classification :

$$y \in \{1, 2, 3, \dots, K\}$$

Many SVM packages already have built-in multi-class classification functionality
 Otherwise, use one-vs-all method

Train K SVM's one to distinguish $y = i$ from the rest, for $i = 1, 2, \dots, K$) get

$$y=1 \rightarrow \theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$$

Pick class i with largest $(\theta^{(i)})^T x$

⇒ Logistic Regression V/s SVM's :

n = number of features ($x \in \mathbb{R}^{n+1}$); m = number of training examples

↳ If n is large (relative to m) (e.g. $n = m$, $n = 10,000$, $m = 1000$)

Use logistic regression or SVM without a kernel ("linear kernel")

↳ If n is small, m is intermediate (e.g. $n = 1-1000$, $m = 10-1000$)

Use SVM with gaussian kernel

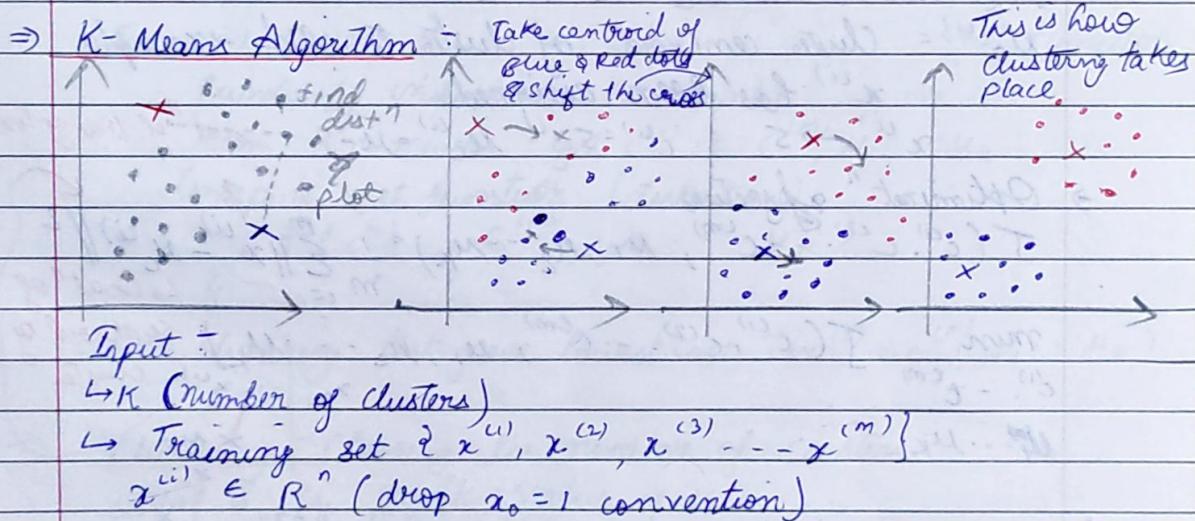
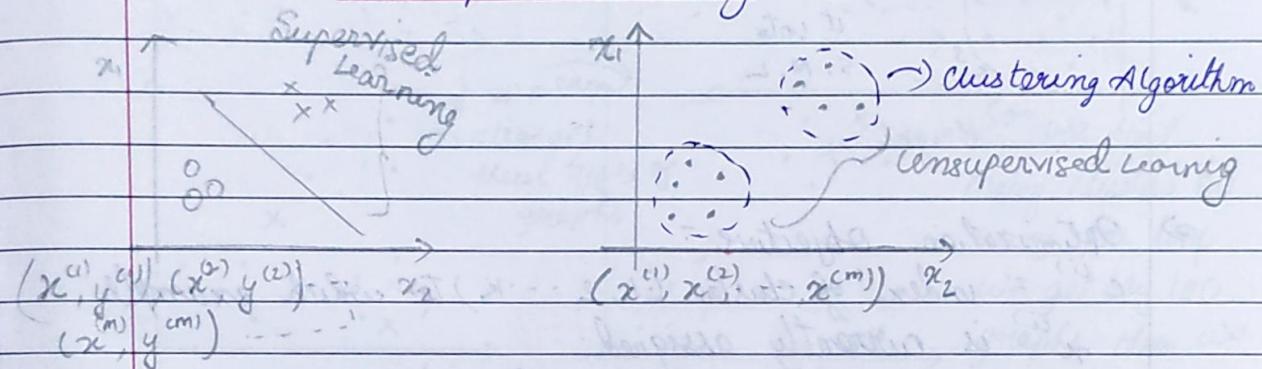
↳ If n is small, m is large (e.g. $n = 1-1000$, $m = 5,000+$)

Create / add more features; then use logistic regression or SVM without a kernel.

↳ Neural networks likely to work well for most of these settings but may be slower to train.

Week - 8

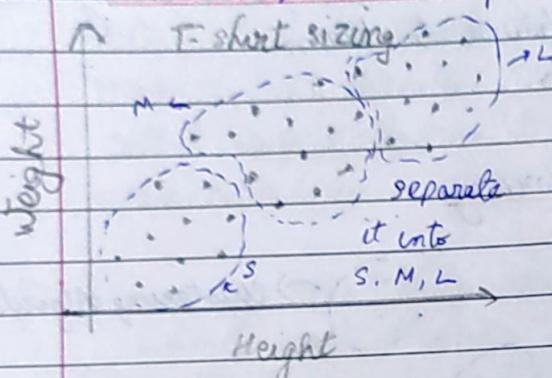
Unsupervised learning



Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$
Repeat {

cluster assignment { for $i=1$ to m minimize $J(\cdot, \cdot)$ w.r.t $\mu_1, \mu_2, \dots, \mu_K$ (keeping $c^{(i)}$ fixed)
closer to red or blue $c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$ $\min_{k} \|x^{(i)} - \mu_k\|^2$
move centroid { for $k=1$ to K moves centroid step
 $\mu_k :=$ average (mean) of points assigned to cluster k
if $\mu_1 = \frac{x_1^{(1)} + x_1^{(5)} + x_1^{(6)} + x_1^{(10)}}{4}$, $\mu_2 = \frac{\text{minimize } J(\cdot, \cdot)}{\text{w.r.t } c^{(1)}, c^{(2)}, \dots, c^{(10)}}$
 $c^{(1)} = 2; c^{(5)} = 2; c^{(6)} = 2; c^{(10)} = 2$

⇒ K-Means for non-separated clusters



⇒ Optimization objective :

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

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

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

$$x^{(i)} \rightarrow S \quad \therefore c^{(i)} = 5 \quad \mu_{c^{(i)}} = \mu_5 \quad \text{locat^n of 5th row}$$

⇒ Optimization objective :

$$J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{c^{(1)}, \dots, c^{(m)}} J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_k)$$

\$x_1\$ \$x_2\$ \$x_3\$ \$x_4\$ \$x_5\$

↳ locat^n of centroid of cluster

⇒ Random initialization : $\mu_1 = x^{(1)}$ $\mu_2 = x^{(2)}$ \dots

↳ Should have $K < m$

↳ Randomly pick K training

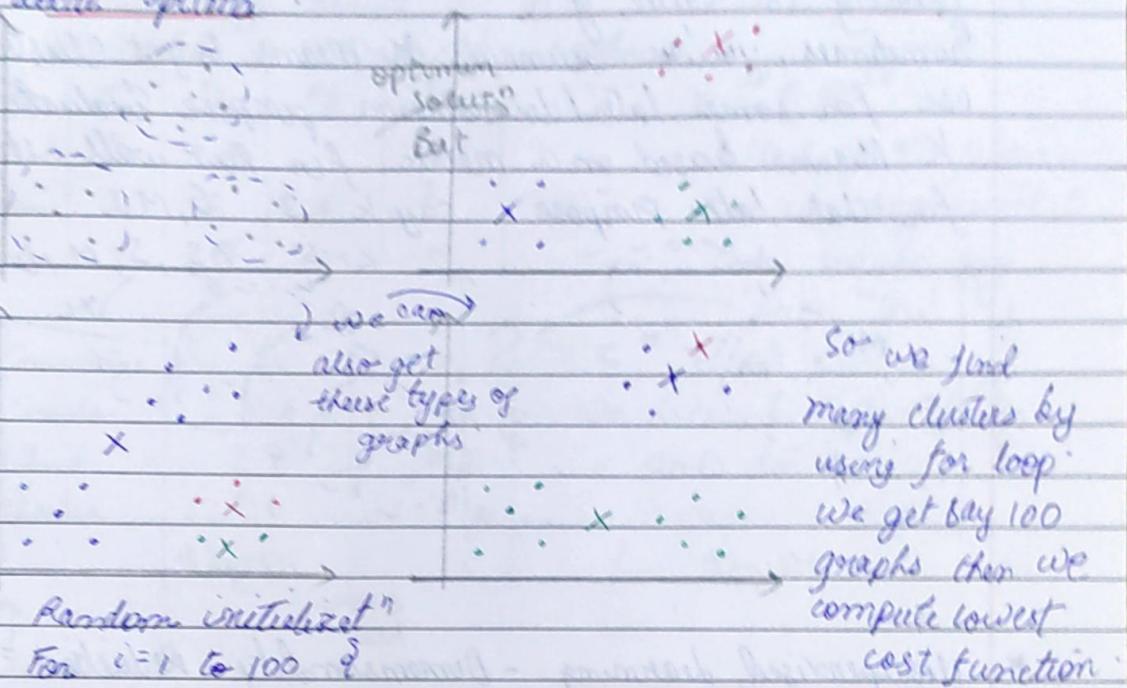
Examples :

↳ Set μ_1, \dots, μ_k equal of these K examples

$K=2$

$X X$

⇒ local optima



Random initialize

For $i=1$ to 100 {

Randomly initialize K-means

Run K-means Get $c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_k$

Compute cost function (distortion)

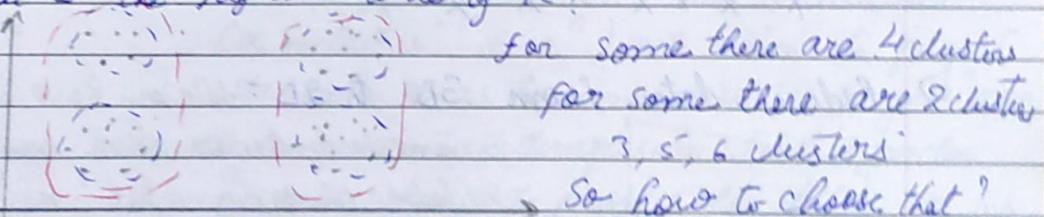
$J(c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_k)$

}

Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, u_1, \dots, u_k)$

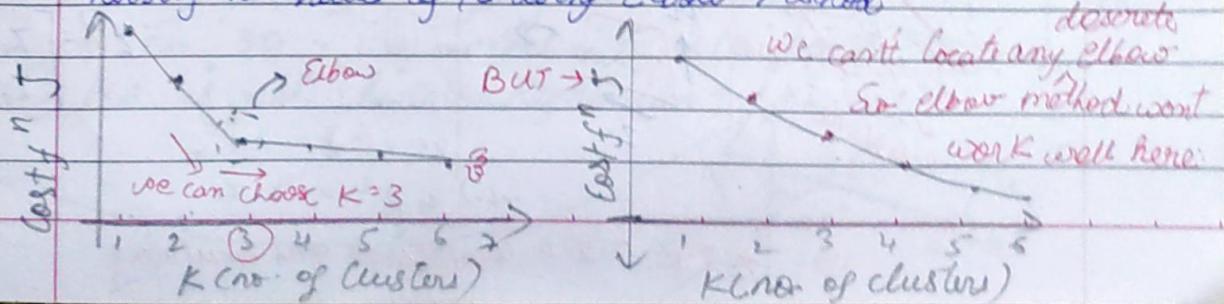
⇒ Clustering = Choosing the number of clusters

What is the right value of K?



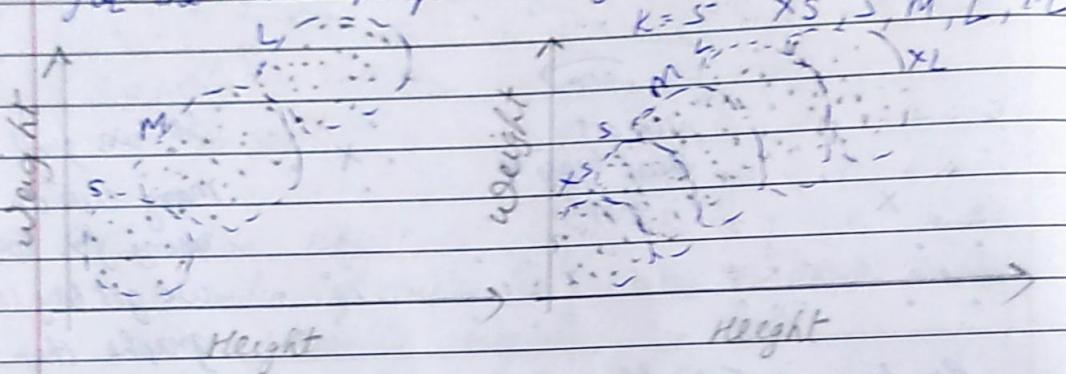
So how to choose that?

⇒ Choosing the value of K using Elbow Method

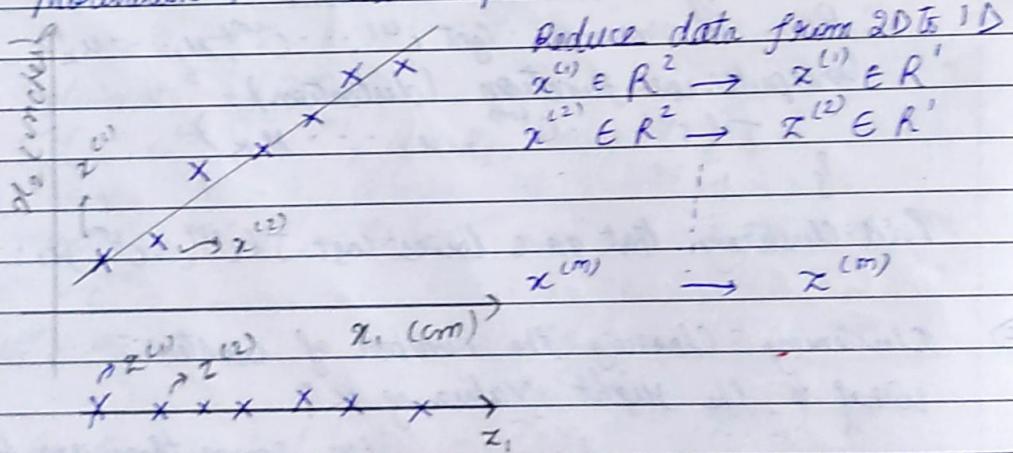


⇒ Choosing the value of K :

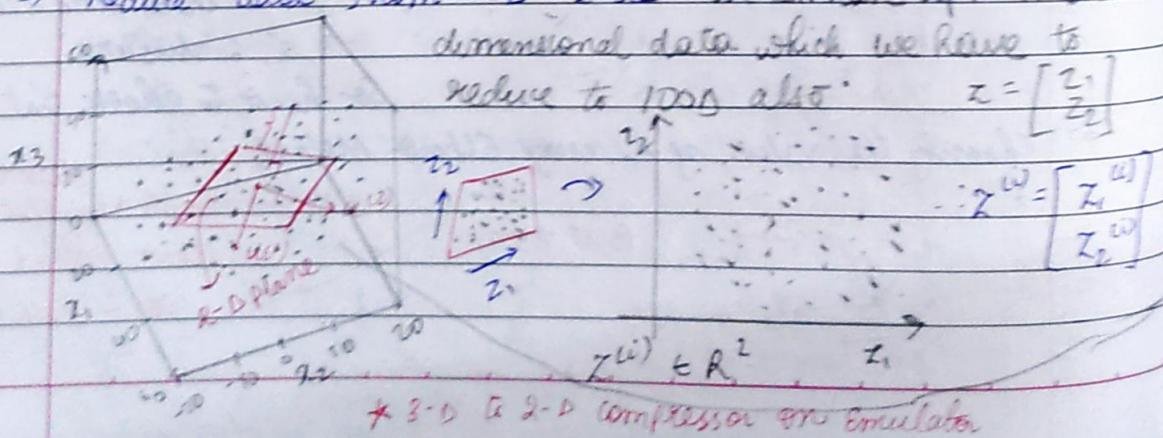
Sometimes, you are running K-Means to get clusters to use for some later/demonstration purpose. Evaluate K-Means based on a metric for how well it performs for that later purpose: say k = 3 S, M, L



* Unsupervised learning - Dimensionality Reduction:
Motivation 1 = Data compression



⇒ Reduce data from 3D to 2D = We can have upto 1000D



⇒ Dimensionality Reductⁿ

Motivation 2 : Data Visualization

Country	x_1	x_2	x_3	x_4
Canada	2	3	4	7
China	3	2	5	8
India	4	6	9	10

$$x^{(i)} \in \mathbb{R}^{50}$$

Instead of 50 features we can represent it two maybe say like this

Country	x_1	x_2
Canada	1.6	1.2
China	1.7	0.3
India	1.6	0.2

$$x^{(i)} \in \mathbb{R}^2$$

We Reduced data from 50D to 2D

x_2
say per person
economy

x_1 , say GDP

⇒ Principal Component Analysis Problem formulation = (PCA)

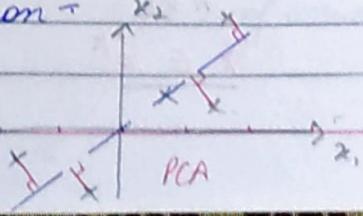
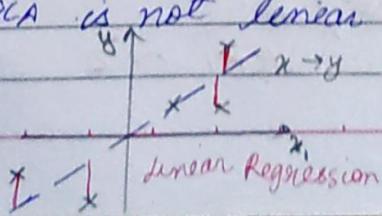
$x \in \mathbb{R}^2$

Case 1: Acc to PCA we will select case 1 as the square distance from the line is much more smaller than Case 2. while line $u^{(1)}$ Reduce from 2D to 1D - Find a direction ($\text{a vector } u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projected error.

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

→ For 3D, we convert it to 2D so $k=2$.

→ PCA is not linear Regression



⇒ PCA Algorithm:

Before using Algorithm we must always do Data preprocessing

↳ Data preprocessing:

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

Preprocessing (feature scaling / mean normalization)

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{Mean of each feature}$$

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

If different features on different scales (e.g. x_1 = size, x_2 = bedrooms), scale features to have comparable range of values. We do feature scaling

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

$$s_j$$

⇒ Algorithm:

Reduce data from n-dimensions to k-dimensions

represents Compute "covariance matrix":

$$\text{matrix} \rightarrow \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \quad \begin{array}{l} \text{product is } n \times n \text{ matrix} \\ \text{greek alpha letter sigma} \end{array}$$

compute "eigen vectors" of matrix Σ :

$$[U, S, V] = \text{svd}(\text{Sigma}); \quad (\text{SVD} \rightarrow \text{singular value decomposition})$$

We can also use $\text{eig}(\text{sigma})$

$$U = \begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(n)} \end{bmatrix} \quad U \in \mathbb{R}^{n \times n}$$

$$x \in \mathbb{R}^n \rightarrow x \in \mathbb{R}^k \quad \begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(K)} \end{bmatrix} \quad \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}_{n \times K}$$

$$\therefore Z = \begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(K)} \end{bmatrix}^T X = \begin{bmatrix} \dots & u^{(1)} & \dots \end{bmatrix} \quad X \in \mathbb{R}^{n \times n}$$

$$X \in \mathbb{R}^K \quad \begin{bmatrix} & & & \\ & & & \\ & & & \end{bmatrix}_{K \times 1} \quad \text{Value}$$

$$K \times 1 \rightarrow z$$

⇒ PCA Algorithm Summary:

After mean normalization (ensure every feature has zero mean) and optionally feature scaling.

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \quad \rightarrow \text{In octane}$$

$$\text{sigma} = (\sqrt{m}) * x' * x;$$

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

$$U_{\text{reduce}} = U(:, 1:k); \quad \left. \begin{array}{l} \text{grab first} \\ k \text{ columns} \end{array} \right\}$$

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

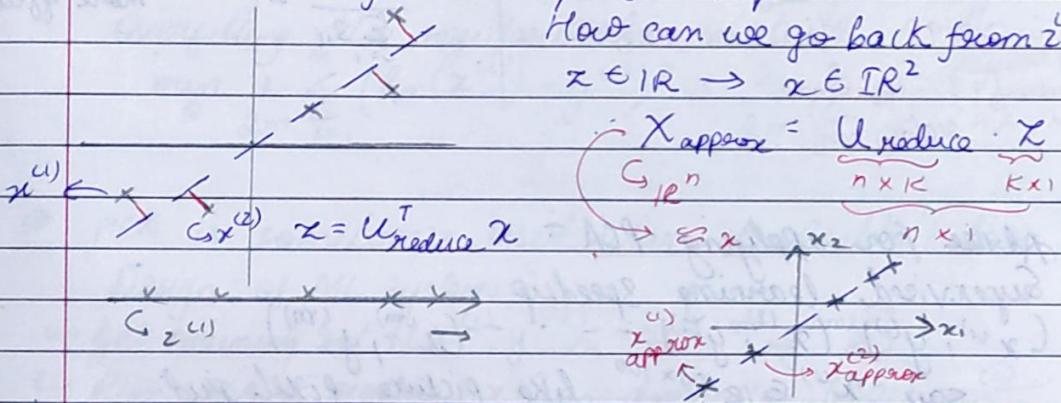
$$* x = \begin{bmatrix} \dots & x^{(1)} & \dots \\ \dots & \vdots & \dots \\ \dots & x^{(m)} & \dots \end{bmatrix}$$

$$\rightarrow z \in \mathbb{R}^n \quad \cancel{x \in \mathbb{R}^n}$$

⇒ Reconstruction from Compressed Representation:

How can we go back from $z \rightarrow x^{(i)}$

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



⇒ Choosing number of principal components (k):

choosing $k \leq$ no. of principal components

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

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

Typically choose k to be smallest value so that

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

$$\leq 0.01 / (1\% / 5\%)$$

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

99 % of variance is retained
95 %

\Rightarrow Choosing k (number of principal components)

Algorithm:

Try PCA with $K = 1/2/3 \dots$ say $K=1$

Compute U_{reduce} , $z^{(1)}, z^{(2)} \dots z^{(m)}$ we

$x^{(1)} \text{ approx} \dots x^{(m)}$

$x^{(1)} \text{ approx} \dots x^{(m)}$

99% of variance

retained

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

$$S = \begin{bmatrix} S_{11} & 0 & \cdots & 0 \\ 0 & S_{22} & 0 & \cdots \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & S_{33} & \cdots \\ \vdots & \vdots & \vdots & \ddots & S_{nn} \end{bmatrix} \quad \text{for } k=2$$

Check if

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

First pick the

smallest value of k for which

$\sum_i S_{ii} \geq 0.99$ then increase K

$\sum_{i=1}^k S_{ii}$ [99% variance retained]

for given K

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \leq 0.01$$

$$\text{or } \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99 \quad \text{if it is much more efficient}$$

\Rightarrow Advice for applying PCA:

Supervised learning Speedup

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

say $x^{(i)} \in \mathbb{R}^{100 \times 100}$ like picture pixel grid

Exact inputs:

Unlabeled dataset: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10000}$

$$x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10000 \text{ (say)}}$$

X
PCA
Z

New training set:

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

Now we can use regression, neural networks... whatever we want.

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined

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

by running PCA only on training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test set.

⇒ Application of PCA :

- ↪ compression : i) Reduce memory/disk needed to store data
ii) Speed up learning algorithm

Choose K by % of variance retain

↪ Visualization : $K = 2$ or 3 as we can only visualize that

⇒ Bad use of PCA : To prevent overfitting : Bad !

Use $\mathbf{z}^{(i)}$ instead of $\mathbf{x}^{(i)}$ to reduce the number of features to $K < n$

Thus fewer features, less likely to overfit.

↪ This might work ok, but isn't a good way to address overfitting. Use regularization instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

regularization is better than PCA for overfitting

⇒ PCA is sometimes used where it shouldn't be :

Design of ML System :

↪ Get training set $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

↪ Run PCA to reduce $\mathbf{x}^{(i)}$ in dimension to get $\mathbf{z}^{(i)}$

↪ Train logistic regression on $\{(\mathbf{z}^{(1)}, y^{(1)}), \dots, (\mathbf{z}^{(m)}, y^{(m)})\}$

↪ Test on test set : Map $\mathbf{x}_{\text{test}}^{(i)}$ to $\mathbf{z}_{\text{test}}^{(i)}$. Run $h_{\theta}(\mathbf{z})$ on $\{(\mathbf{z}_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}), \dots, (\mathbf{z}_{\text{test}}^{(m)}, y_{\text{test}}^{(m)})\}$

BUT ; how about doing the whole thing without PCA ?

Before implementing PCA ; first try running whatever you want to do with original/raw data $\mathbf{x}^{(i)}$. Only if that doesn't do what you want, then implement PCA and consider using $\mathbf{z}^{(i)}$.

Week - 9

*

Anomaly Detection

→ Anomaly detection Example: For aircraft engine features

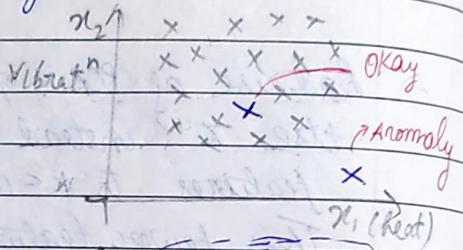
x_1 = heat generated

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

x_2 = vibration intensity

New engine : x_{test}

...



→ Density estimation :

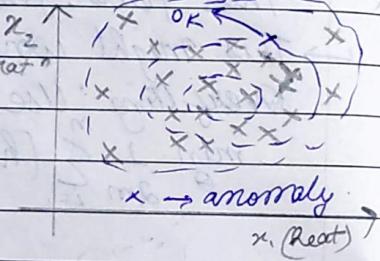
Dataset : $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

Is x_{test} anomalous?

Model $p(x)$

$p(x_{test}) < \epsilon \rightarrow$ flag anomaly

$p(x_{test}) \geq \epsilon \rightarrow$ OK



→ Anomaly detection example :

i) Fraud detect :-

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

↳ Model $p(x)$ from data.

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

ii) Manufacturing

Monitoring computers on a data center

$x^{(i)}$ = features of machine i

x_1 = memory use, x_2 = number of disk accesses/sec

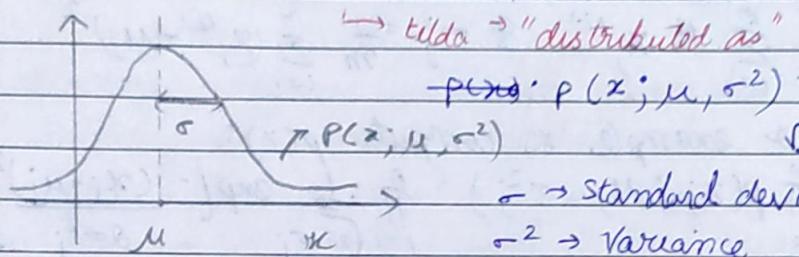
x_3 = CPU load, x_4 = CPU load / network traffic

...
 $p(x) < \epsilon$

\Rightarrow Gaussian Distribution = (Normal)

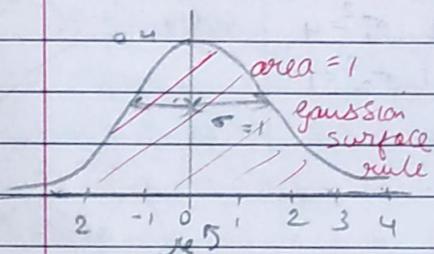
Say $x \in \mathbb{R}$, if x is a distributed Gaussian with mean μ & variance σ^2 . \rightarrow script N stands for Normal = Gaussian

$$x \sim N(\mu, \sigma^2)$$

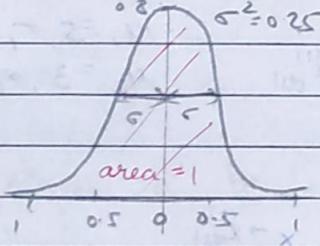


\Rightarrow Gaussian Distribut" Example:

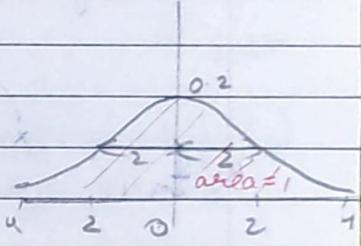
$$\mu = 0, \sigma = 1$$



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



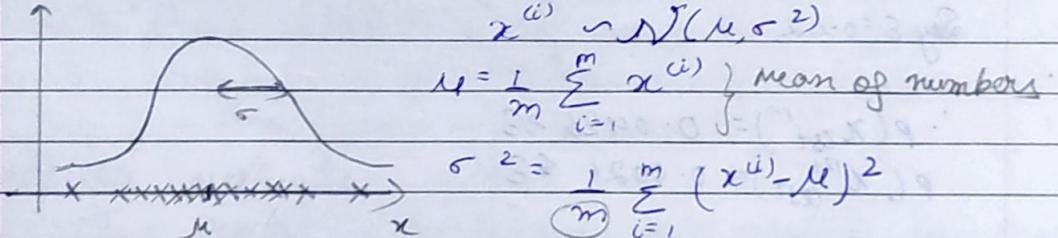
$$\mu = 0, \sigma = 2$$



\Rightarrow Parameter estimation:

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

$$x^{(i)} \sim N(\mu, \sigma^2)$$



$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \text{Mean of numbers}$$

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

In statistics we use $(m-1)$ But in ML we ignore this factor as m is very large

\Rightarrow Algorithm:

Density Estimation - $x \in \mathbb{R}^n$

$$x_i \sim N(\mu_i, \sigma_i^2)$$

Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \quad x_1 \sim N(\mu_1, \sigma_1^2)$

$$p(x) = p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \dots p(x_n; \mu_n, \sigma_n^2)$$

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

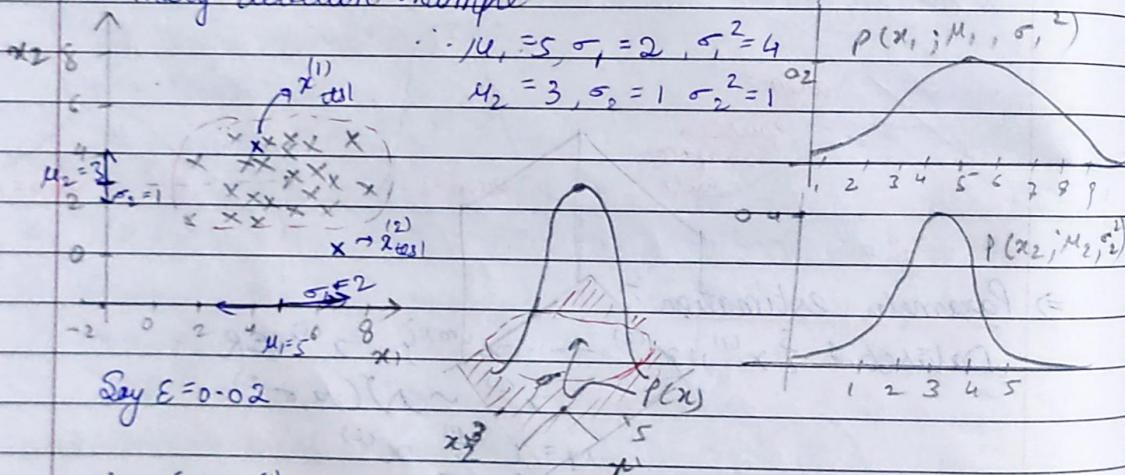
Product Symbol

\Rightarrow Anomaly detection Algorithm:

- Choose features x_i that you think might be indicative of anomalous examples: $\{x^{(1)}, \dots, x^{(m)}\}$
 - Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$
 - Given new example x , compute $p(x)$:

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

\Rightarrow Anomaly detection Example:



$$\therefore p(x_{\text{test}}^{(1)}) = 0.0426 \geq \epsilon$$

$$p(x_{\text{test}}^{(2)}) = 0.0021 \leq \epsilon$$

\Rightarrow Developing and Evaluating an Anomaly Detection System:

\hookrightarrow Importance of real number evaluation

When developing a learning algorithm (choosing features etc) making decisions is much easier if we have a way of evaluating our learning algorithm.

\hookrightarrow Assume we have some labeled data, of anomalous and non-anomalous examples ($y=0$ if normal, $y=1$ if anomalous).

- ↳ Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples / not anomalous)
- ↳ cross-validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}) \dots (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
Test set: $(x_{test}^{(1)}, y_{test}^{(1)}) \dots (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$
- ⇒ Aircraft engines monitoring Examples:
10000 good (normal) engines
20 flawed engines (anomalous) 2-50
- ↳ Training set: 6800 good engines ($y=0$) $P(x) = P(x; \mu_1, \sigma_1^2)$
CV: 2000 good engines ($y=0$), 10 anomalous ($y=1$)
Test: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Alternative - ~~X~~

- ~~X~~ Training set: 6000 good engines
- CV & Test: 4000 good engines ($y=0$), 10 anomalous ($y=1$)
Some put same 4K data in CV & Test which is not good

⇒ Algorithm Evaluation:

- ↳ Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$
- ↳ On a cross-validation / test example x , predict:

$$y = \begin{cases} 1 & \text{if } p(x) \leq \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases}$$

$y = 1$ if $p(x) \leq \epsilon$ (anomaly)
 $y = 0$ if $p(x) \geq \epsilon$ (normal)

↳ Possible evaluation metrics:

True positive, false positive, false negative, true negative
Precision / recall

F-Score

Can also use cross-validation set to choose parameter ϵ .
Select diff. value of ϵ to get maximum F-Score.

⇒ Anomaly detection v/s supervised learning :-

Anomaly detection

- Very small number of positive → large number of positive & examples ($y=1$) : ($0-20$ is common) negative examples
- Large number of negative ($y=0$)
- Examples: PCB like aircraft
- Many different "types" of anomalies → enough positive examples for algorithm - Hard for any algorithm - other to get a sense of what to learn from positive examples positive examples are like, what the anomalies look like; future positive examples likely future anomalies may look to be similar to ones in nothing like any of the anomalies training set.
- examples we've seen so far:

↳ Fraud detection

↳ Email spam classification

↳ Manufacturing (e.g. aircraft engines)

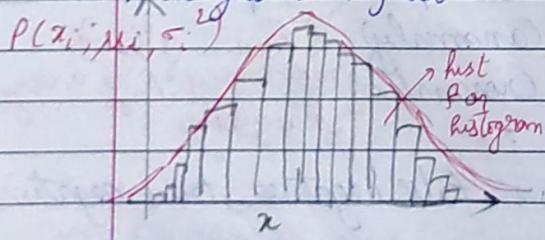
↳ Cancer classification

↳ Monitoring machines in a data center

↳ Weather prediction (summer, rain, winter)

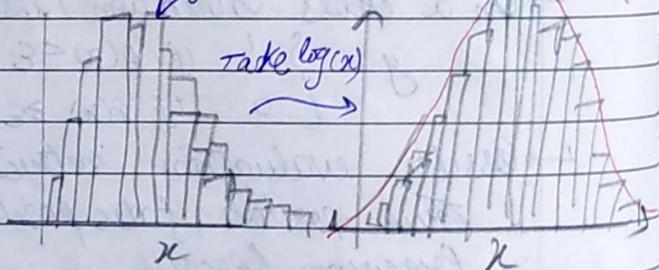
⇒ Choosing what features to use:-

Non-gaussian feature



Non gaussian curve

This looks much more gaussian



$x_i \rightarrow \log(x_i)$ / $x_2 \rightarrow \log(x_2 + 1) / \log(x_2 + c) \dots$ can be used

like $x_3 \rightarrow \sqrt{x_3} / x_3^k$ / $x_4 \rightarrow x_4^{1/k}$ To look data gaussian

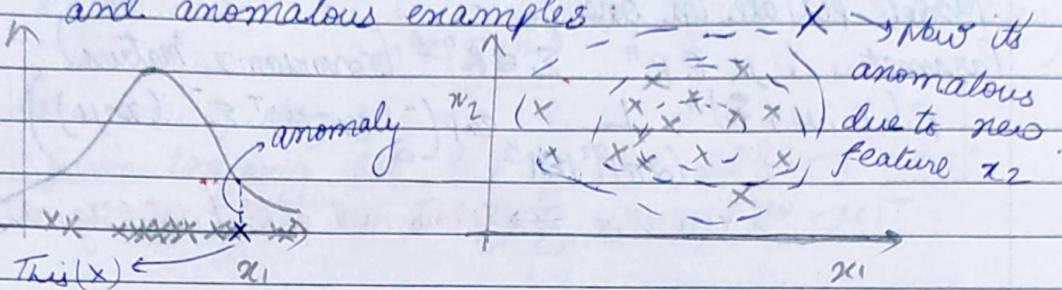
⇒ Error analysis for anomaly detection :-

↳ Want $p(x)$ large for normal example x .

$p(x)$ small for anomalous example x .

→ Most common problem :-

$p(x)$ is comparable (say, both large) for normal and anomalous examples :-



Show be normal but we find it anomaly So we can add new features to it say (x_2) & make a new plot & then make that plot in use & show that the a new example is shown anomalous

⇒ Monitoring computers in a data center :-

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

x_1 = memory use of computer.

x_2 = number of disc accesses/sec

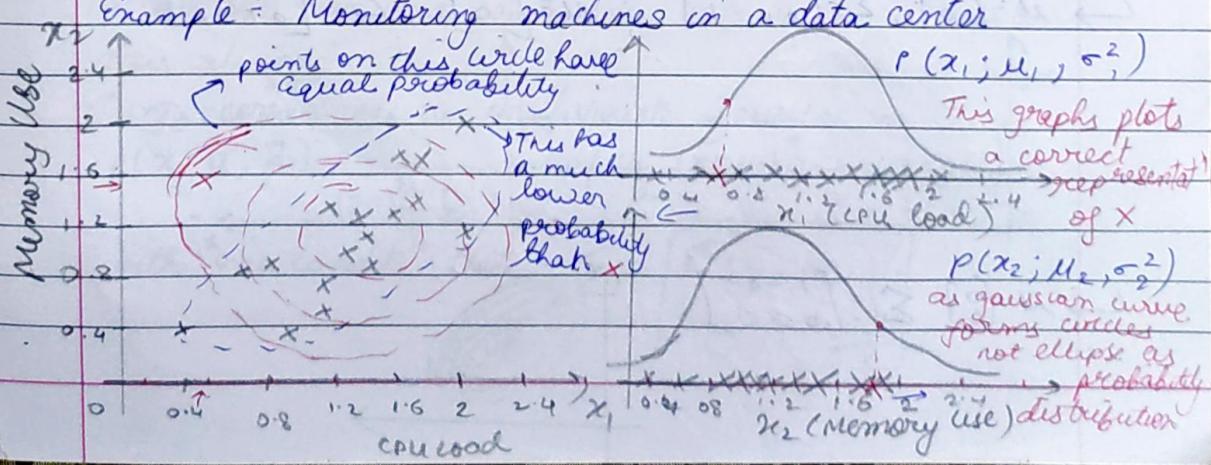
x_3 = CPU load

x_4 = network traffic } If it gives a bigger value or or being dependent on each other

New x_5 = $\frac{\text{CPU load}}{\text{network traffic}}$ / $x_6 = (\text{CPU load})^2$

⇒ Multivariate Gaussian Distribution :-

Example : Monitoring machines in a data center



⇒ Multivariate Gaussian (Normal) distribution :

$x \in \mathbb{R}^n$ Don't model $p(x_1), p(x_2), \dots$ etc. separately.

Model $p(x)$ all in one go :

Parameters : $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix).

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

(determinant of Σ / $\det(\Sigma)$)

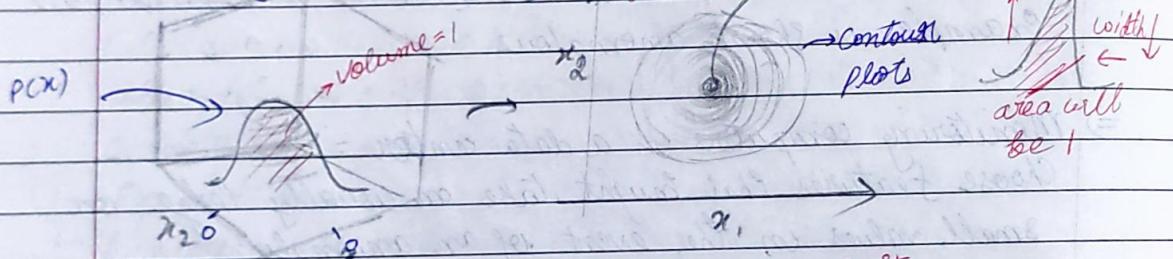
⇒ Multivariate Gaussian (Normal) Examples :

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

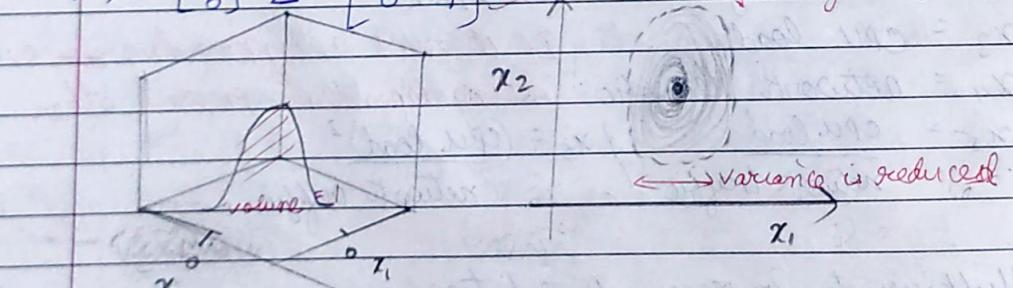
Center peak can be altered for $\Sigma = \begin{bmatrix} 0.6 & 0 & 0 \\ 0 & 0.6 & 0 \\ 0 & 0 & 0.6 \end{bmatrix}$

By changing μ peak height ↑

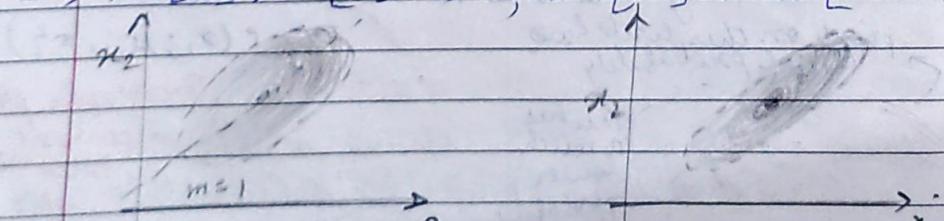
height ↑



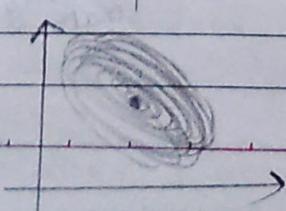
$$\hookrightarrow \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix} \rightarrow \text{reduces variance of 1st feature } x_1, \text{ variance of 2nd feature } x_2$$



$$\hookrightarrow \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}; \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



$$\hookrightarrow \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$



⇒ Anomaly Detection using the Multivariate Gaussian Distribution

Parameters μ, Σ ; $x \in \mathbb{R}^n$ & $\Sigma \in \mathbb{R}^{n \times n}$

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

Parameter fitting :-

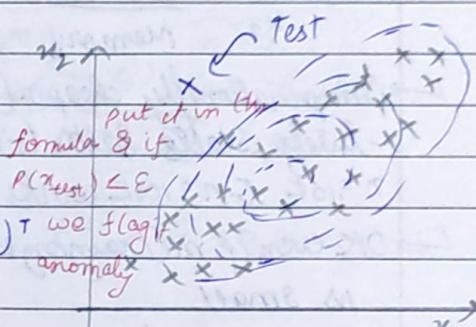
Given training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \rightarrow \mathbb{R}^n$

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

i) Fit model $p(x)$ by setting

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

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



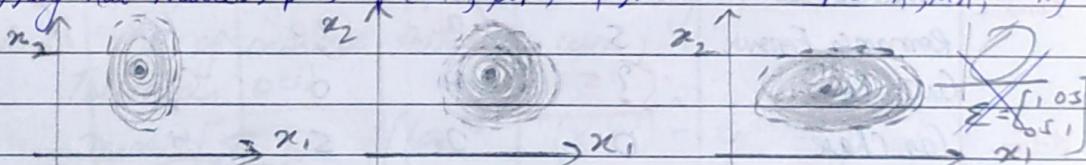
ii) Given a new example x , compute

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

Flag an anomaly if $p(x) < \epsilon$.

⇒ Relationship to Original model :-

Original model; $p(x) = p(x_1; \mu_1, \sigma_1^2) x_1 \cdots x_n p(x_n; \mu_n, \sigma_n^2)$



The contours are along same axis as x_1 & x_2 . These are axis aligned.

∴ It corresponds to multivariate Gaussian.

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

with one condition that

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_n^2 \end{bmatrix}$$

Original Model

$$P(x_i | \mu_i, \Sigma_i) \propto P(x_i; \mu_i, \Sigma_i) \rightarrow P(x_i; \mu_i, \Sigma_i) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu_i)^T \Sigma_i^{-1} (x_i - \mu_i)\right)$$

↳ Manually creates features to capture anomalies where x_1, x_2 correlations between take unusual component "of value".

$$\begin{matrix} x_3 = x_1 \\ x_2 = \text{CPU load} \\ \quad \quad \quad \text{memory} \end{matrix}$$

$\Sigma \in \mathbb{R}^{n \times n}$ & for big n Σ^{-1} expensive

↳ Computationally cheaper (alternatively, scale better to larger n) Σ singular

$$n=10 \text{ & } m=100 \text{ works fine}$$

$$\Sigma \sim 2^2 \text{ parameters } / m_1 = x_2 / m_2 = x_3 = x_4 + x_5$$

↳ OK even if m (training set size) is small

↳ Must have $m > n$ or else Σ is non-invertible or $m \leq 10$

* Recommender System

Problem formulation

Example: Predict movie ratings

User rates movies using one - five stars.

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)	
Romantic	5	5	0	0	$n_u = 4$
Love at Last	5	(?) 4.5	(?) 0	0	$n_m = 5$
Romance Forever	5	4	0	(?) 0	
Cute Love	(?) 5	4	0	(?) 0	
Action	0	0	5	4	
Car Chase	0	0	5	4	
Swords V/s Karate	0	0	5	(?) 4	

$$n_u = \text{no. of users}$$

$$n_m = \text{no. of movies}$$

$$x(i, j) = 1 \text{ if user } i \text{ has rated movie } j$$

$$y^{(i, j)} = \text{rating given by user } i \text{ to movie } j$$

$$0 \dots 5 \quad (\text{defined only if } x(i, j) = 1)$$

We have to predict the ?

⇒ Content Based Recommendation System :

Same table as left

$$\begin{array}{l} \text{---} \\ x_1 = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix} \end{array}$$

Movie	$\theta^{(1)}$ Alice(1)	$\theta^{(2)}$ Bob(2)	$\theta^{(3)}$ Carol(3)	$\theta^{(4)}$ Dave(4)	$Z_1(\text{correct})$	$Z_2(\text{act}^n)$
1 Love at last	5	5	0	0	0.9	0.3
2 Romance forever	5	?	?	0	1.0	0.01
3 Cuti Love	24.95	4	0	?	0.99	0
4 Car Chase	0	0	5	4	0.1	1.0
5 Swords vs Karate	0	0	5	?	0	0.9

$$\Theta^{(1)} \in \mathbb{R}^{n+1} \quad n=2$$

For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user j 's rating for movie i with $(\theta^{(j)})^T x^{(i)}$ stars.

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \quad \Theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad \therefore (\Theta^{(1)})^T x^{(3)} = 5 \times 0.99 = 4.95$$

(say)

\Rightarrow Problem formulation:

$g(i, j) = 1$ if user j has rated movie i (0 otherwise)

$y_{(i,j)}^{(c(i,j))}$ = rating by user j on movie i (if defined)

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

$x^{(i)}$ = feature vector for movie i

For user j , movie i , predicted rating $(\theta^{(i)})^T(x^{(j)})$

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

To learn $\Theta(j)$:

$$\min_{\Theta^{(j)}} \left[\sum_{i: x(i,j)=1} \left((\Theta^{(j)})^T (x^{(i)}) - y^{(i,j)} \right)^2 \right] \frac{1}{2\alpha^{(j)}} + \frac{\lambda}{2\alpha^{(j)}} \sum_{k=1}^n (\Theta_k^{(j)})^2$$

~~Q⁽ⁱ⁾ SR~~ \rightarrow To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n)}$.

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

$\underbrace{\quad \quad \quad}_{J(\theta^{(1)}, \dots, \theta^{(n_u)})}$

$$\Rightarrow \text{Gradient descent update} \div \\ \theta_c^{(j)} := \theta_c^{(j)} - \alpha \sum_{i: n(i,j) \neq 1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k=0)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i: n(i,j) \neq 1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) / x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

$\frac{\partial}{\partial \theta_k^{(j)}} \sum_{i: n(i,j) \neq 1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) + \lambda \theta_k^{(j)}$

\Rightarrow Collaborative Filtering \div

Problem motivation: What if we don't know x_1 & x_2 ? $x_0 = ?$

movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)	Romance	Action
$x^{(1)}$ Love at last	5	5	0	0	?	0.0
$x^{(2)}$ Romance forever	5	?	?	0	?	?
$x^{(3)}$ Cute love	?	4	0	?	?	?
$x^{(4)}$ Car chase	0	0	5	4	?	?
$x^{(5)}$ Speedy instant	0	0	5	?	?	?

$$\theta_1 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad \theta_2 = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad \theta_3 = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} \quad \theta_4 = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

so that $(\theta^{(1)})^T x^{(1)} \approx 5$ & $(\theta^{(3)})^T x^{(4)} \approx 0 \quad \left\{ \begin{array}{l} x^{(1)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ x^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \end{array} \right.$

\Rightarrow Optimization Algorithm \div

Given $\theta^{(1)}, \dots, \theta^{(n_m)}$ to learn $x^{(1)}, \dots, x^{(n_m)}$:

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

Given $\theta^{(1)}, \dots, \theta^{(n_m)}$ to learn $x^{(1)}, \dots, x^{(n_m)}$:

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

⇒ According to collaborative filtering :-

Given $x^{(1)} \dots x^{(n_m)}$ (and movie ratings),

we can estimate $\theta^{(1)} \dots \theta^{(n_u)}$

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

We can estimate $x^{(1)} \dots x^{(n_m)}$

Guess $\theta \rightarrow x \rightarrow \theta \rightarrow x \dots \theta \in R^n$

$x_0 = 1 \quad x \in R^n$ as x_0 is cancelled

⇒ Collaborative Filtering optimization objective :-

Given $x^{(1)} \dots x^{(n_m)}$ estimate $\theta^{(1)} \dots \theta^{(n_u)}$:

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

Given $\theta^{(1)} \dots \theta^{(n_u)}$ estimate $x^{(1)} \dots x^{(n_m)}$:

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

Minimizing $x^{(1)} \dots x^{(n_m)}$ and $\theta^{(1)} \dots \theta^{(n_u)}$ simultaneously :

$$J(x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)})^2$$

$$+ \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_m} \sum_{k=1}^n (\theta_k^{(i)})^2$$

$$\min_{x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)}} J(x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)})$$

$$\theta^{(1)} \dots \theta^{(n_u)}$$

⇒ Collaborative Filtering algorithm :-

- Initialize $x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)}$ to small random values.
- Minimize $J(x^{(1)} \dots x^{(n_m)}, \theta^{(1)} \dots \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm).

Eg for every $j=1 \dots n_u, i, l=1 \dots n_m$;

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

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

There is no $x_0 \& \theta_0$.

(iii) For a user with parameters Θ & a movie with (learned) features x ; predict a star rating of $\theta^{(i)}$.

→ Vectorization : low-Rank Matrix Factorization.
 What can we do with this algorithm? a given one product can we find other products related to it?

Collaborative filtering

$$Y = \begin{bmatrix} 5 & 8 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & 0 \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \quad \text{Predicted Rating} = (l_{ij}) \rightarrow (\Theta^{(i)})^T (x^{(j)}) \\ = (\Theta^{(n_u)})^T (x^{(n_m)})$$

$$X = \begin{bmatrix} - (x^{(1)})^T \\ \vdots \\ - (x^{(n_m)})^T \end{bmatrix} \quad \Theta = \begin{bmatrix} - (\Theta^{(1)})^T \\ \vdots \\ - (\Theta^{(n_u)})^T \end{bmatrix}$$

∴ we can compute
 predicted rating $= X \cdot \Theta$
 low-rank matrix factorization.

→ Finding related movies

For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$
 $x_1 = \text{Romance}; x_2 = \text{Action}; x_3 = \text{Comedy}; x_4 = \dots$

→ How to find movies j related to movie i ?

Small $\|x^{(i)} - x^{(j)}\|$ → movies i & j are similar

→ 5 most similar movies to movie i :

Find the 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$

⇒ Implementational Detail: Mean Normalization

User who have not rated any movies -

Movie	A(1)	B(2)	C(3)	D(4)	Eve(5)	
Love at last	5	5	0	0	? 0.25	$\begin{bmatrix} 5 & 5 & 0 & 0 & ? \end{bmatrix}$
Romance forever	5	?	?	0	? 2.5	$\begin{bmatrix} 5 & ? & ? & 0 & ? \end{bmatrix}$
Cute love	?	4	0	?	? 0.2	$\begin{bmatrix} ? & 4 & 0 & ? & ? \end{bmatrix}$
Car Chase	0	0	5	4	? 2.25	$\begin{bmatrix} 0 & 0 & 5 & 4 & ? \end{bmatrix}$
Swords v/s Karate	0	0	5	?	? 0.125	$\begin{bmatrix} 0 & 0 & 5 & 0 & ? \end{bmatrix}$

$$\min_{\theta^{(5)} \in \mathbb{R}^2} \frac{1}{2} \sum_{(i,j) : r(i,j) = 1} ((\theta^{(5)})^T x^{(i)})^2 + \frac{\lambda}{2} \sum_{i=1}^n \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^m \sum_{k=1}^n (\theta_k^{(j)})^2$$

No note for Eve as Eve has not rated any movie

for $n=2$; $\theta^{(5)} \in \mathbb{R}^2$ for Eve(5)

$$\frac{1}{2} [(\theta_1^{(5)})^2 + (\theta_2^{(5)})^2]$$

If we minimize it we will get $\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

∴ We must do $(\theta^{(5)})^T x^{(i)} = 0$

So we will rate is 0 everywhere. X Not good

⇒ Mean normalization :-

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}, M = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & 1.25 & ? \end{bmatrix}$$

Now we will use this

Y as and learn it $\theta^{(j)}$ & $x^{(i)}$

For user j on movie i predict

$$(\theta^{(j)})^T (x^{(i)}) + \mu_i$$

$$\text{(User 5 (Eve))} : (\theta^{(5)})^T (x^{(i)}) + \mu_i \quad \therefore \text{We can say that}$$

$$\theta^{(5)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

he will predict the average of all values

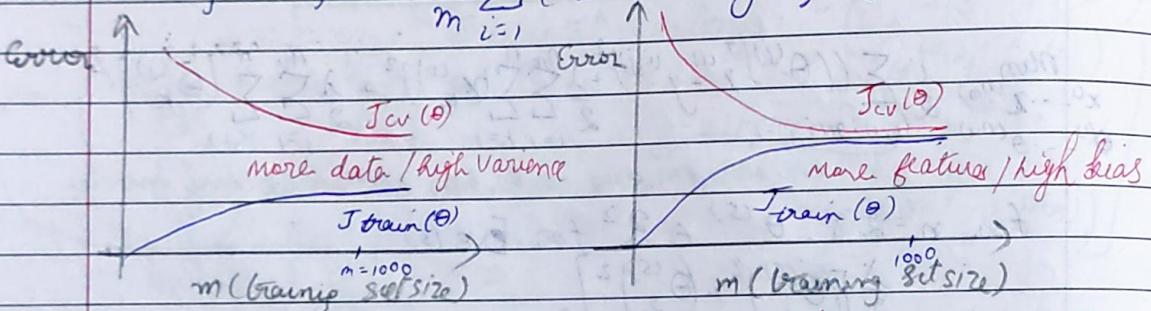
Week - 10

- * Gradient Descent with large Datasets:
 - "It's not who has best algorithm that wins
 - "It's who has the most data"

\Rightarrow Learning with large Datasets:

$$m = 100,000,000 \text{ (new example)}$$

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



If we have millions of data we can plot learning curves to check whether more data is required or more features are required.

\Rightarrow Stochastic Gradient Descent:

this is a modification of gradient descent for big m

$$h_\theta(x) = \sum_{j=0}^n \theta_j x_j$$

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

Repeat ?

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

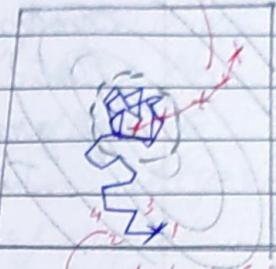
(for every $j = 0, \dots, n$) { If m is large this will take time
say $m = 300,000,000$

\rightarrow Batch Gradient Descent

Stochastic Gradient Descent

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

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



i) Randomly shuffle dataset say (1 - 100)

ii) Repeat for $i = 1 \dots m$

$$\theta_j := \theta_j - \alpha (\hat{y}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

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

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

Instead of summing up in batch GD here we are iterating & improving with each step. Not accurate but gives a value surrounding it.

Mini-Batch Gradient Descent

Batch GD: use all m examples in each iteration

Stochastic GD: use 1 example in each iteration

Mini-batch GD: Use b examples in each iteration

b = mini-batch size

b can be 10 - 100

$$m=1000; \text{ get } b=10 \text{ examples } (x^{(i)}, y^{(i)}) \dots (x^{(i+9)}, y^{(i+9)})$$

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (\hat{y}(x^{(k)}) - y^{(k)}) \cdot x_j^{(k)}$$

$i := i + 10$

Mini-batch gradient descent

Say $b = 10$; $m = 1000$

Repeat for

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

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (\hat{y}(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

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

⇒ Stochastic Gradient Descent convergence :-

Checking for convergence

Plot $J_{\text{train}}(\theta)$ as a function of the number of iterations of gradient descent.

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

$$m = 300,000,000$$

⇒ Stochastic GD :-

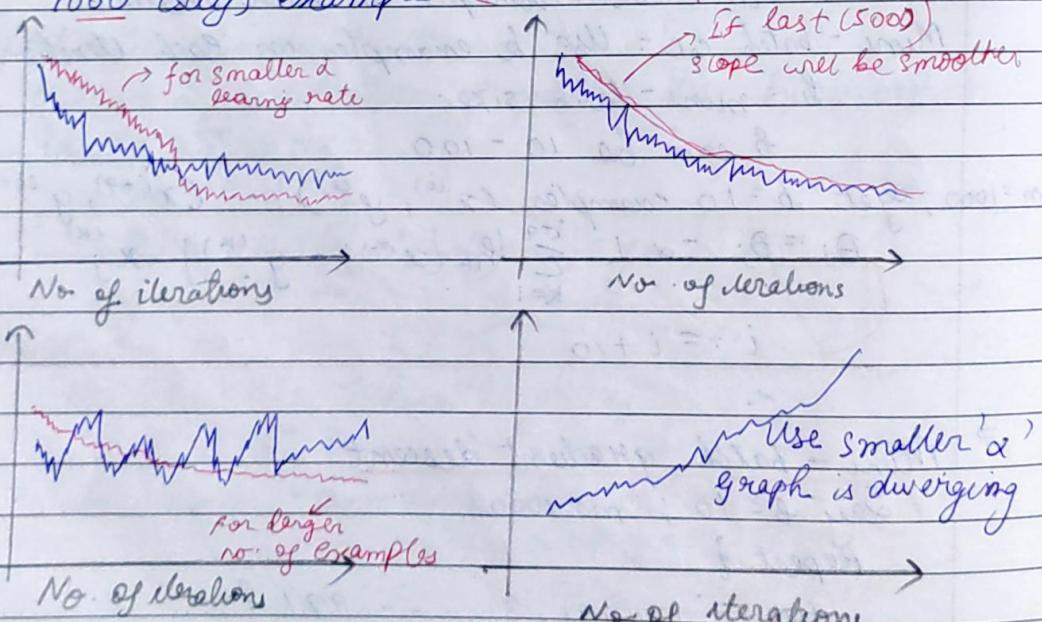
$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

During learning compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$.

Every 1000 iteration (say) plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algo.

⇒ Checking for convergence :-

Plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$, averaged over the last 1000 (say) example



Learning rate α is typically held constant. Can slowly decrease α over time if we want θ to converge [Eg. $\alpha = \text{const.}$]

as iteratⁿ numbr ↑ α so it gets close to global minima

iteratⁿ numbr + const 2

\Rightarrow Online Learning :

- \hookrightarrow Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price and users sometimes choose to use your shipping service ($y=1$), sometimes not ($y=0$).
 - \hookrightarrow Features x capture properties of user of origin/destination and asking price. We want to learn $p(y=1|x; \theta)$ to optimize price.
- \uparrow price we can use logistic regression

Repeat forever {

we use this update it

get (x, y) corresponding to user.

update θ using (x, y)

$$\theta_j := \theta_j - \alpha (h_\theta(x) - y) x_j \quad (j = 0 \dots n)$$

\uparrow Can adapt to changing user preferences or else update won't be good.

\Rightarrow Online learning Example :

Product search (learning to search)

User searches for "Android phone 1080 p camera".

Have 100 phones in store. We will return 10 results.

- \hookrightarrow $y=1$ if user clicks on link ; $y=0$ otherwise
- learn $p(y=1|x; \theta)$ } predicted click through rate (CTR)

\hookrightarrow Use to show user the 10 phones they're most likely to click on.

Other examples : Choosing special offers to show user, customized selection of news articles ; product recommendation ...

\Rightarrow Map Reduce and Data Parallelism:

Batch gradient descent: $m = 400$ this can be applied to 4×10^9

$$\theta_j := \theta_j - \alpha \sum_{i=1}^{400} (\text{logistic}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Machine 1: Use $(x^{(1)}, y^{(1)}) \dots (x^{(400)}, y^{(400)})$

$$\left[\begin{array}{c} x_j^{(i)} \\ y^{(i)} \end{array} \right] \quad \text{temp}_j^{(1)} = \sum_{i=1}^{100} (\text{logistic}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

$$\left[\begin{array}{c} x_j^{(i)} \\ y^{(i)} \end{array} \right] \quad \text{Machine 2: Use } (x^{(101)}, y^{(101)}) \dots (x^{(200)}, y^{(200)})$$

$$\text{temp}_j^{(2)} = \sum_{i=101}^{200} (\text{logistic}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Say $m=400$ ($x^{(m)}, y^{(m)}$)

Same for machine 3 & 4.

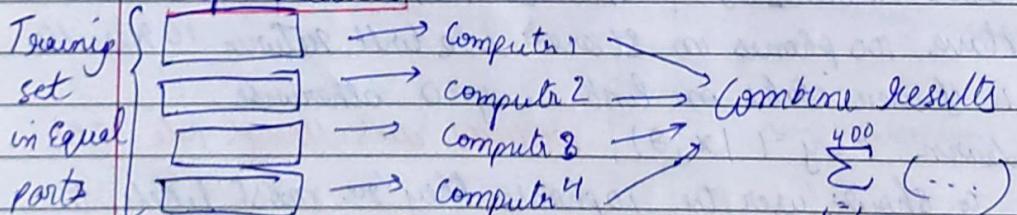
$$\text{temp}_j^{(1)}, \text{temp}_j^{(2)}, \text{temp}_j^{(3)}, \text{temp}_j^{(4)}$$

combine

$$\theta_j := \theta_j - \alpha \frac{1}{400} (\text{temp}_j^{(1)} + \text{temp}_j^{(2)} + \text{temp}_j^{(3)} + \text{temp}_j^{(4)})$$

i.e. sum of all the parts
for $j = 0, \dots, n$

\Rightarrow Map - Reduce:



\Rightarrow Map - reduce and summation over the training set:

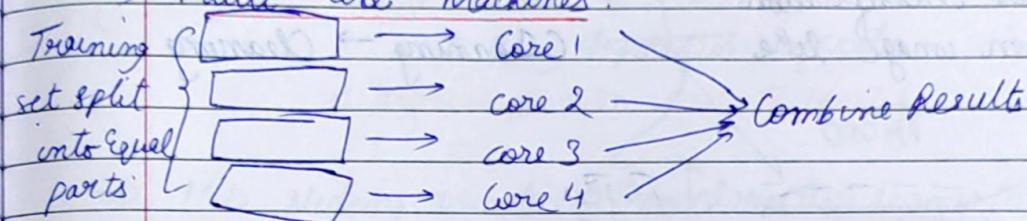
Many learning algorithms can be expressed as computing sums of functions over the training set.

Ex for advanced optimization with logistic regression need:

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

$$\hookrightarrow \frac{\partial J_{\text{train}}(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

⇒ Multi-core Machines:



Week - 11

* Photo OCR = Photo optical character recognition

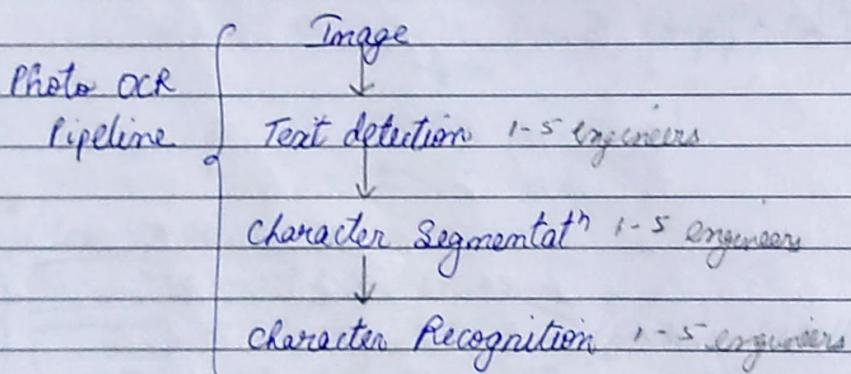
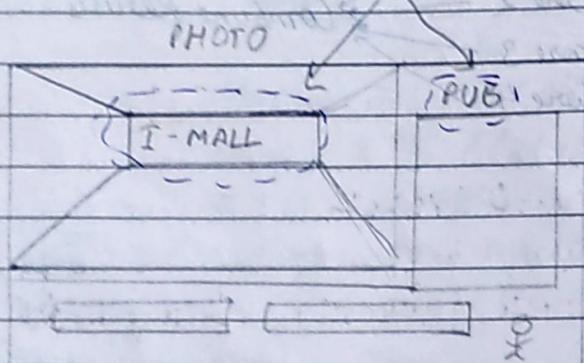
=> Photo OCR Pipeline :-

i) Text detection

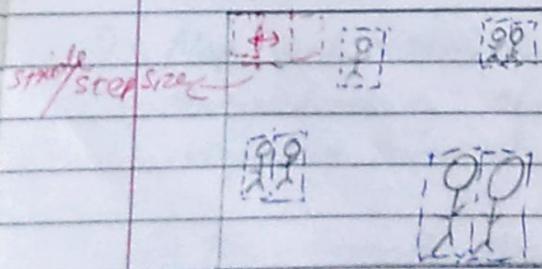
ii) Character Segmentation → PUB T - MALL

iii) Character Classification → P → P U → u B → B

iv) Correction image like → C1eaning → Cleaning



→ Sliding Windows :-



First we give computer a lot of images of pedestrians & non-pedestrians. Then we make a small window and move it over all the photo & try to detect pedestrians.

↳ We can select different size of detecting Pedestrians windows too & it identifies pedestrians.

- We can also perform it for text detection like goes train your computer with lots of text & non-text as training set. Then it identifies the text from the image.

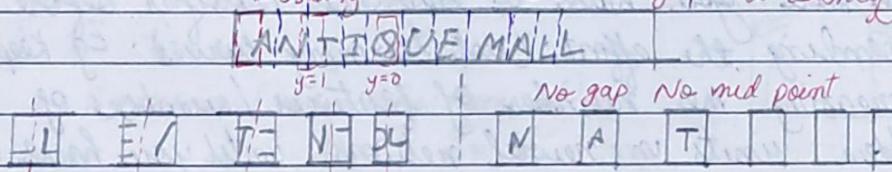
* This is how it is identified

Using Expansion operator

* After expansion easy to identify

- ⇒ 1-D sliding window for character segmentation :-

Sliding this is the gap to identify



Positive Examples ($y=1$)

Negative Examples ($y=0$)

Train it and then test it on Antigie Mail using Sliding window.

- ⇒ Getting lots of Data and Artificial Data :-

↳ Character Recognition :-

* Letter Recognition using different fonts too

↳ Synthesizing by introducing Distortion

* Using distortion to create more training set

↳ Synthesizing data by introduct "distort" :- Speech Recognition

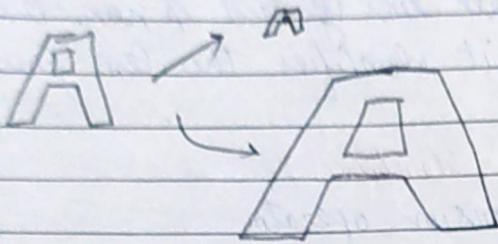
* Speed addit " to create more training set

↳ Distortion introduced should be representation of the type of noise/ distortions in the test set :-

$$\begin{array}{c} A \rightarrow A \\ \downarrow \\ A \end{array}$$

These two can be a good training set, for audio we can add background noise, beeps ...

- Usually does not help to add purely random/ meaningless noise to your data.



x_i = intensity (brightness) of pixel i
 $x_i \leftarrow x_i + \text{random noise}$
 By increasing size, intensity
 won't count as good
 training set

- ⇒ Discussion on getting more data:-

- i) Make sure you have a low bias classifier before expending the effort. (plot learning curves). Eg keep increasing the number of features / numbers of hidden units in neural network until you have a low bias classifier.
- ii) "How much work would it be to get 10x as much data as we currently have?"
- iii) Artificial data synthesis → hours? say $m = 1000$
 - ↳ Collect / label it yourself → 10 sec / Example
 - ↳ "Crowd source" (ex Amazon Mechanical Turk) $m = 10,000$ sec few days may be

- ⇒ Ceiling Analysis : What part of the pipeline to work on next?
 Estimating the errors due to each component (ceiling analysis)

Image

what part of the pipeline
should you spend the

Text Detect'

most time trying to improve?

Character Segmentation

Now we can manually correct the errors

Character Recognition

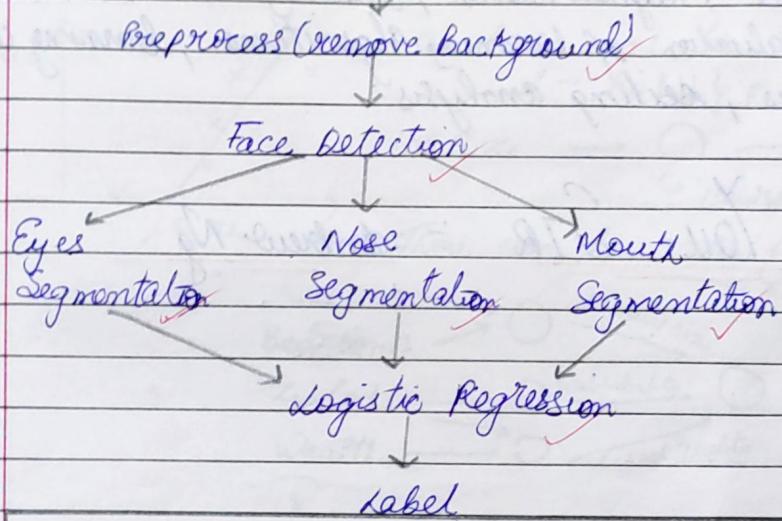
Manually label it

Component	Accuracy for overall system after correcting the pipelines)
Overall System	72%
Text detection	89% ↓ 1%
Character Segmentation	90% ↓ 1%
Character Recognition	100% ↓ 1%

⇒ Another ceiling analysis Example :
Face recognition from images (Artificial Example)

Let's make a pipeline first :

Camera Image



Component	accuracy	notes
Overall System	85%	↓ 0.1% ↗ not a worth thing
Preprocess (Remove backg ⁿ)	85.1%	↓ 5.1% ↗ to spend on time
Face Detection	91%	↓ 4%
Eyes Segmentation	95%	↓ 1%
Nose Segmentat ⁿ	96%	↓ 1%
Mouth Segmentat ⁿ	97%	↓ 3%
Logistic Regression	100%	

* Summary :-

- i) Supervised learning :- $(x^{(i)}, y^{(i)})$
 - ↳ linear Regression, logistic Regression, neural Networks, SVM's
- ii) Unsupervised learning :- $(x^{(i)})$
 - ↳ K-Means, PCA, Anomaly Detection
- iii) Special Applications / Special topics :-
 - ↳ Recommender Systems, Large scale Machine learning
- iv) Advice on building a machine learning system.
 - ↳ Bias / variance, regularization, deciding what to work on next :- evaluation of learning algorithms, learning curves, error analysis, ceiling analysis

THANK YOU SIR :- Andrew Ng