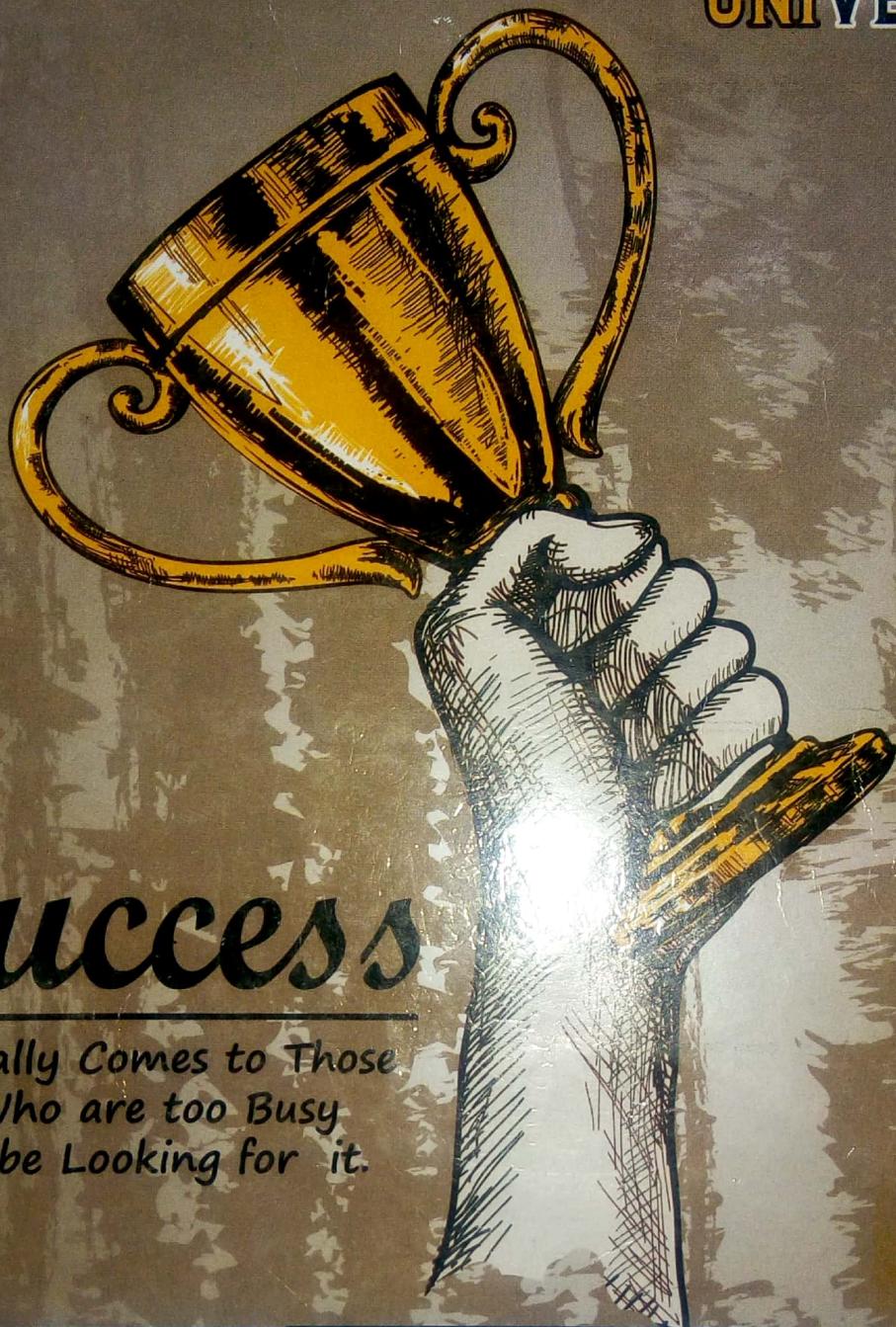


UNIVERSITY



Success

Usually Comes to Those
Who are too Busy
to be Looking for it.

University

Exercise Book
Write Your Mission

Name Naima Islam Nodi

School/College BAUST

Class 4-1 Section B Roll No. 15010

Subject Machine Learning Year 2018

1058

06.05.18.

Udacity → Machine learning

Arthur Samuel:

The field of study that gives computers the ability to learn without being explicitly programmed.

Tom Mitchell:

The field of machine learning is concerned with the quest of how to construct computer program that automatically improve with experience. A com. prog. is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

i.e. checkers Game:

E = the experience of playing

T = the task of playing checkers

P = the probability that the program will win the next game.

Machine learning algorithms:

1. Supervised learning

2. Unsupervised learning → Semi Supervised.

3. Reinforcement learning

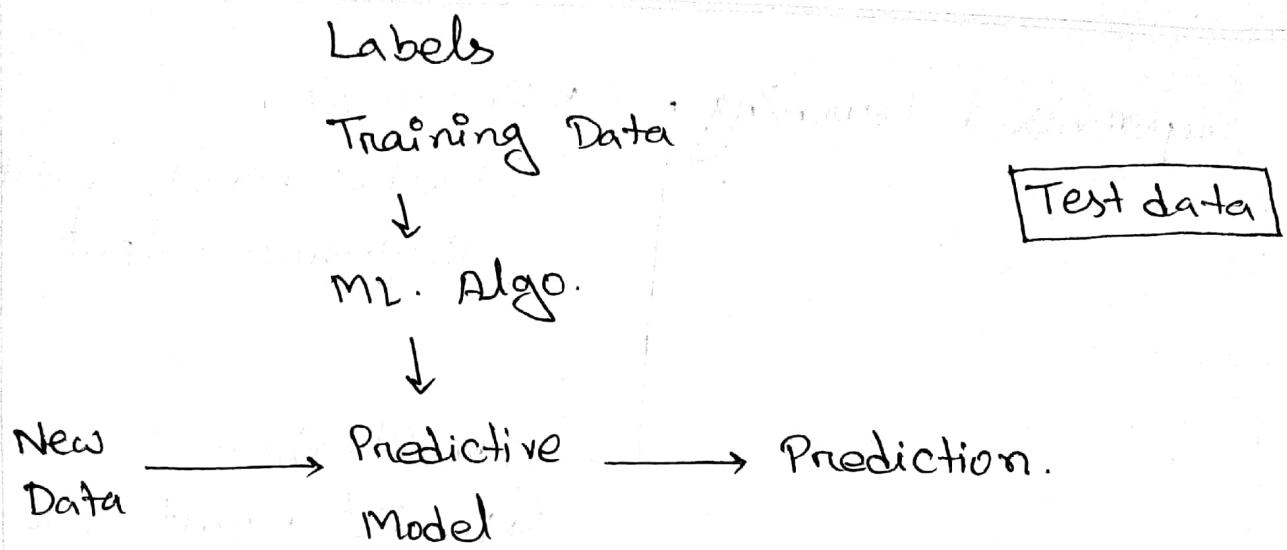
Acervous Vs. Non-Acervous

1. Supervised Learning:

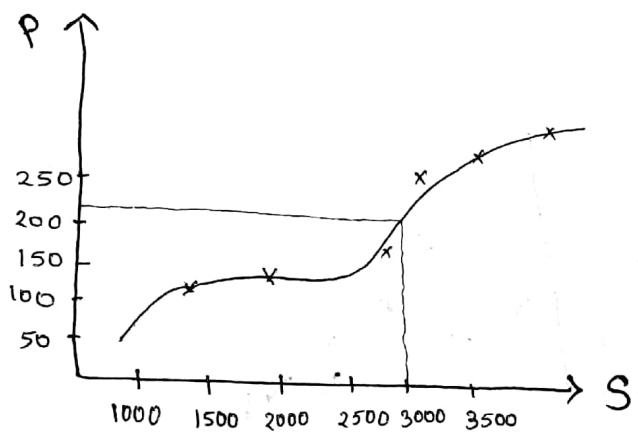
SL is where you have input variables (x) and output variables (y) and you use an algorithm to learn the mapping function from the input to the output.

$$y = f(x) \quad (\text{level data})$$

The goal is to approximate the mapping function so well that when you have new input data (x) that you can predict output variable (y) for that data.



Model

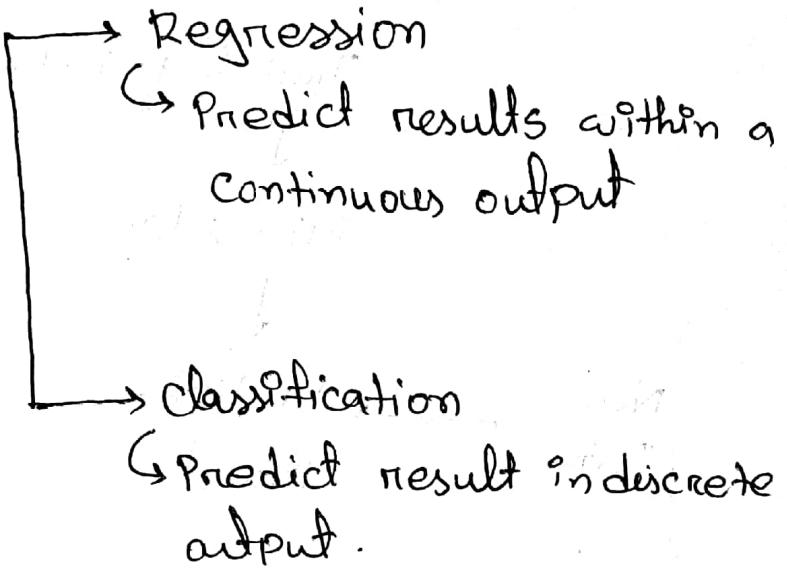


$$Y = f(x)$$

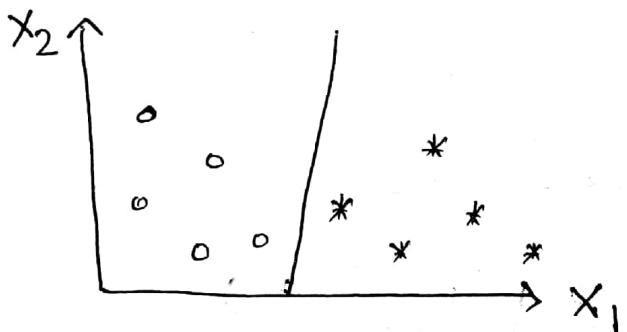
$$P = f(S)$$

08.05.18

Supervised Learning



* Classification:



* Unsupervised Learning:

Here we have little or no idea what our result should look like.

input variable(x) but no Y	Underlying Structure
-------------------------------	----------------------

Unsupervised Learning

→ Clustering

↳ want to discover the inherent grouping in the data such as group customers by purchasing behavior. (UL classification)

→ Association

labeled → supervised

* Clustering:

- Share certain degree of similarity.
- more dissimilar to object in other clusters.

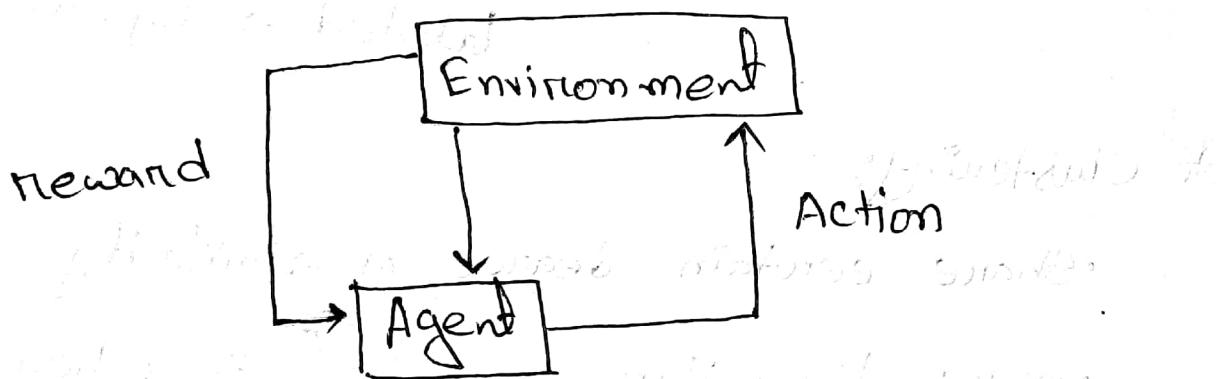
* Association:

↳ Discover rules that describe large portions of our data.

↳ such as : People that buy X also tend to buy Y.

* Reinforcement Learning

→ The goal is to develop a system (agent) that improve its performance based on interactions with the environment, no answer key.



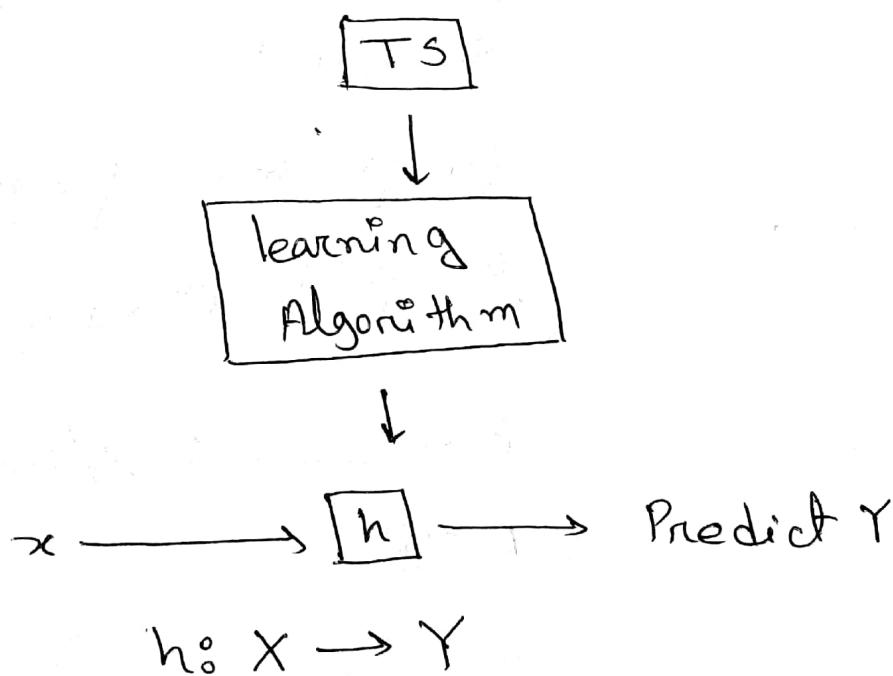
13.05.18

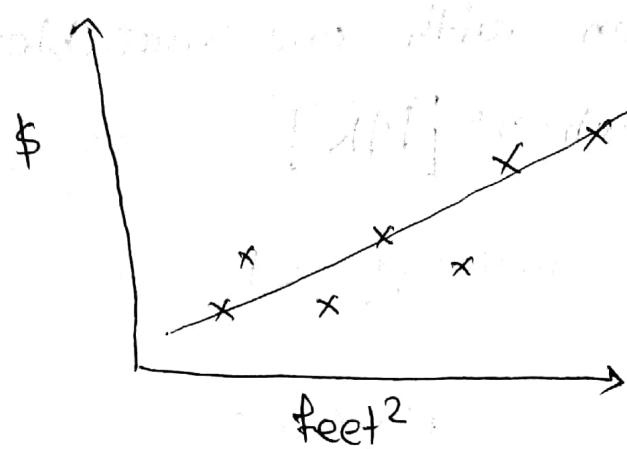
* linear regression with one variable:

Model Representation: [MR]

TS → Training Set

<u>Feet</u> ²	<u>1000 \$s</u>
2104	400
1600	330
2400	369
1416	540
:	:





*Univariate linear regression.

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

m = # No. of training example

n = # No. of variables

x = input

y = output

(x, y) = training example

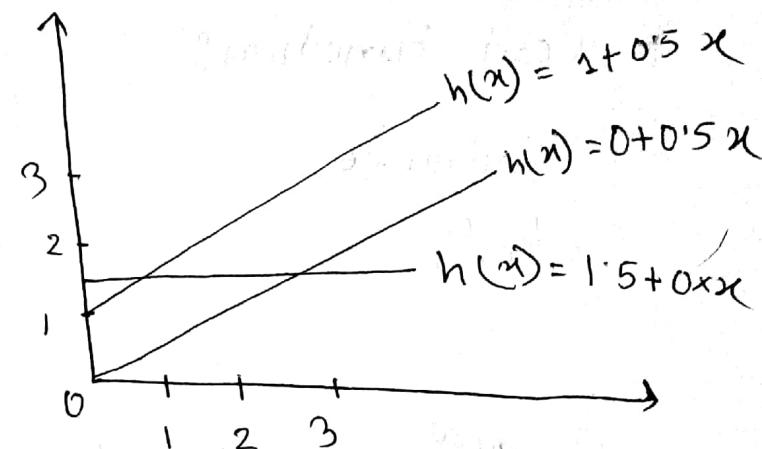
$(x^{(i)}, y^{(i)})$ = i^{th} training example

θ → Parameters/weights

$$\theta \in \mathbb{R}$$

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

$$\begin{array}{c|c} \theta_0 = 1.5 & \theta_0 = 0 \\ \hline \theta_1 = 0 & \theta_1 = 0.5 \end{array}$$

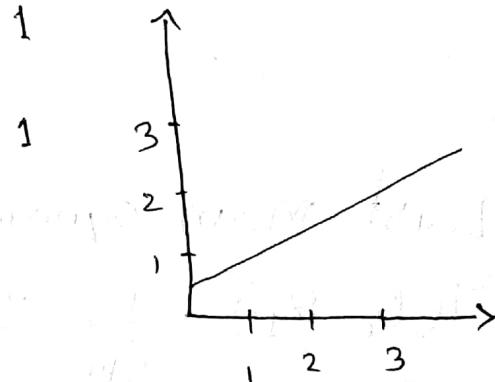


$$\theta_0 = 0 \quad 0.5 \quad 1 \quad 1$$

$$\theta_1 = 1 \quad 1 \quad 0.5 \quad 1$$

i.e.

X	Y
0	4
1	7
2	7
3	8



$$h_{\theta}(x) = \theta_0 + \theta_1 (x) ; \quad \theta_0 = \theta_1 = 2$$

$$= 2 + 2 * x$$

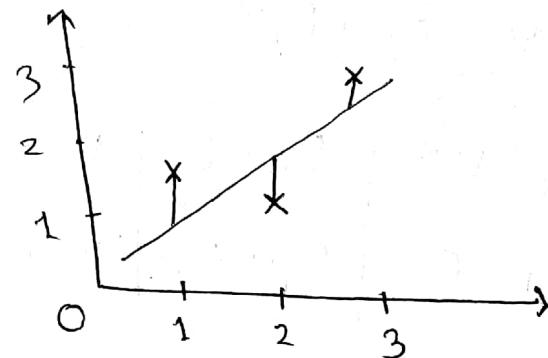
$$\text{if } x = 0 \rightarrow 2 + 2 \cdot 0 = 2$$

$$1 \rightarrow 2 + 2 \cdot 1 = 4$$

* Cost Function:

minimize

θ_0, θ_1



Minimize

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

Least Square Error.

Least Mean Square Error (LMS):

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

Minimize

θ_0, θ_1

$$\text{Hypothesis: } h_{\theta}(x) = \theta_0 + \theta_1 x$$

Parameters: θ_0 and θ_1

Cost function:

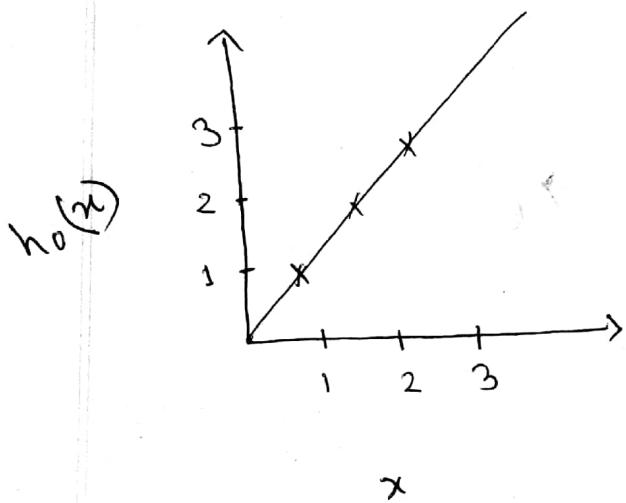
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Goal: minimize

$$\theta_0, \theta_1 \quad J(\theta_0, \theta_1)$$

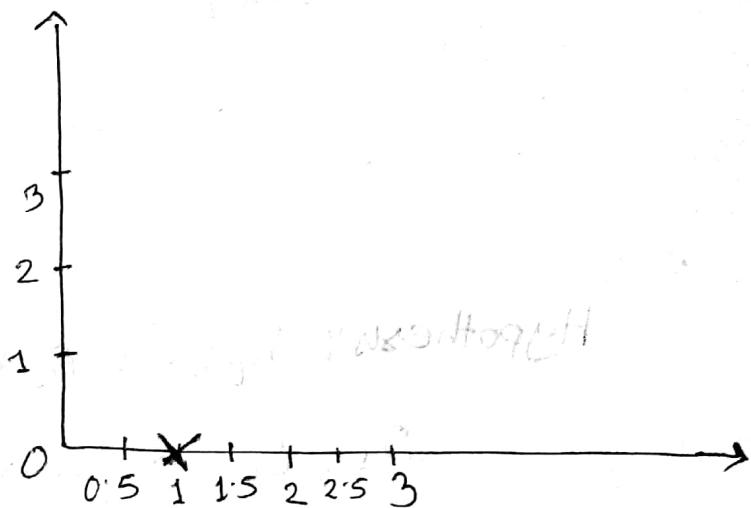
$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \text{Let, } \theta_0 = 0$$

$$h_{\theta}(x) = \theta_1 x$$

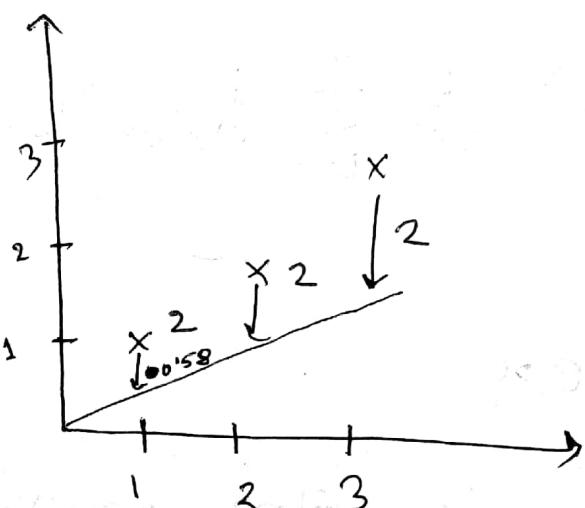


$$\text{if } \theta_1 = 1$$

$$\text{then, } h_{\theta}(x) = x \quad [Y = x]$$

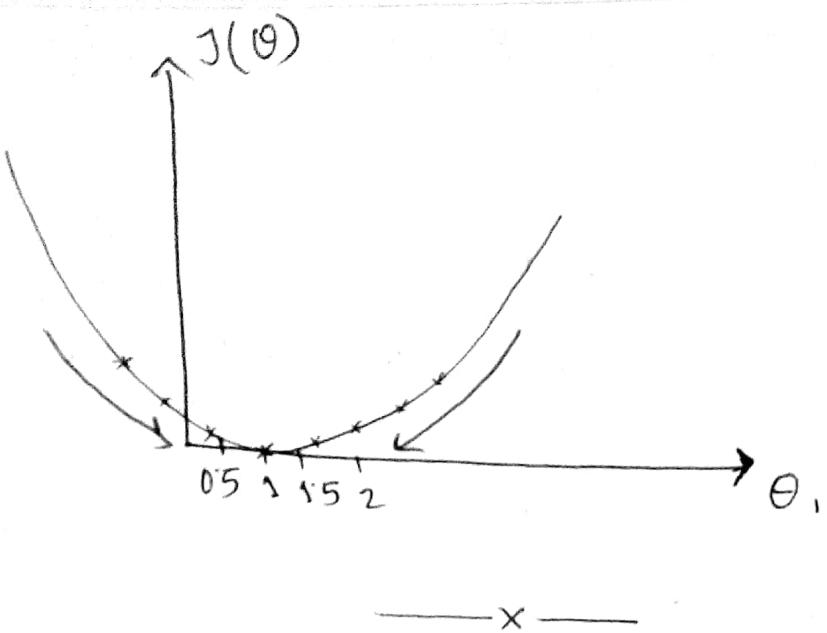


\downarrow position of 1
As distance 0 (1-1=0)



$$\theta_1 = 0.5$$

$$J(0.5) = \frac{1}{2m} \left[(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2 \right] \\ = \frac{1}{2 \times 3} () = 0.58$$



15.05.18

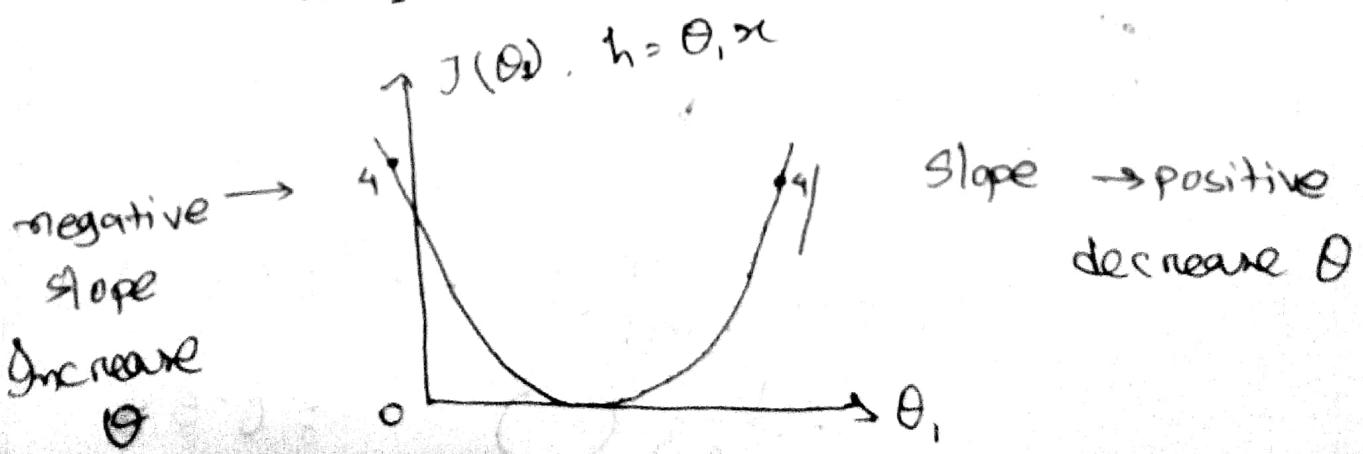
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

$\theta_1, \theta_0 \rightarrow$ Parameters

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

Goal: Minimize

θ_0, θ_1



Gradient Descent:

→ Start with θ_0, θ_1 (any value)

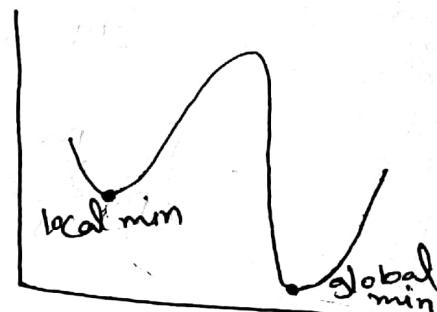
→ Keep changing θ_0, θ_1 , until we have minimum

$J(\theta)$ (Hopefully)

→

If we get more
global solution

repeat {



problem in
gradient descent

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta), \quad j = 0, 1$$

}

repeat

α = learning rate

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

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

}

When $j=0$,

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

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

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

$$= \frac{1}{2m} \cdot 2 \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_0} (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})$$

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

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

When, $j=1$,

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

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

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

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

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

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

repeat

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

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

$$\} \quad \theta_0 - \theta_1 \geq \epsilon \quad \text{or} \quad J(\theta_1) - J(\theta_0) \geq \epsilon \quad \dots \quad // \text{Condition to stop.}$$

update simultaneously:

$$\{ \quad \text{temp0} = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta)$$

$$\text{temp1} = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta)$$

$$\theta_0 = \text{temp0}$$

$$\theta_1 = \text{temp1}$$

}

Batch Gradient descent:

↓ Solution

Incremental / Stochastic Gradient descent:

(first, but possibility of getting global solution is less.)

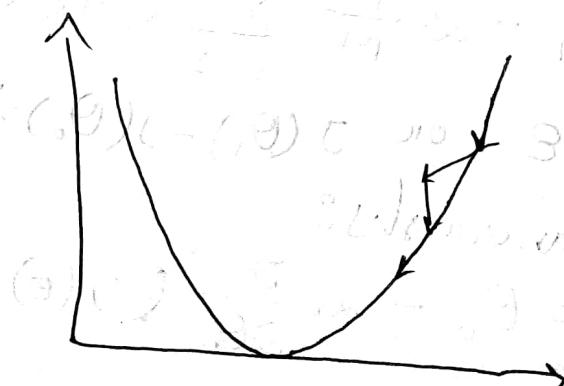
Loop {

for i = 1 to m {

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

}

; for j = 0, 1



— X — Optimal theta

"Special class"

17.05.18.

eliademy.com

~~Python~~

Introduction to Machine Learning and Decision Tree

aktarhossain@dafodilvarsity.edu.bd

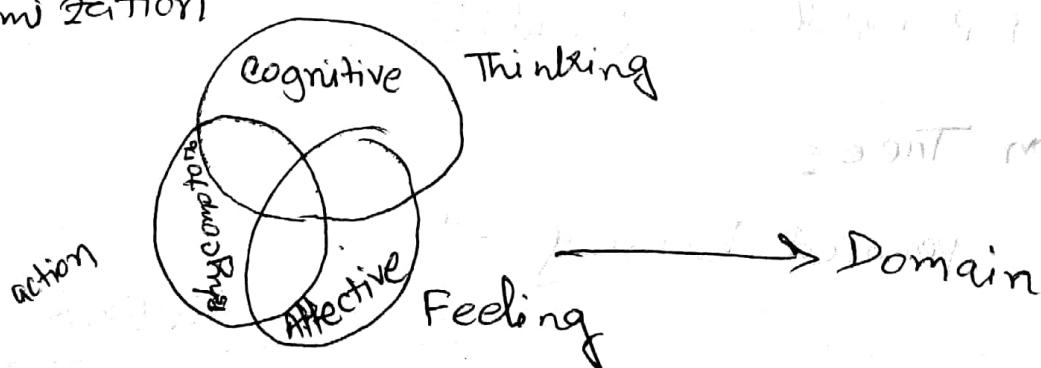
Industry 4.0

Three components of ML

Representation

Evaluation

Optimization



Bloom's Taxonomy of learning

Supervised or inductive learning:

markany.com

backward search \rightarrow Power
 2^{100} combinations
every single type is tuple

Python Basics.pdf.

Everything in Python is an object.

1. get-pip.py

Palindrome

Install pip installer

C:\> Python get-pip.py (Internet)

2. Install Numpy & sci-pi

C:\> pip install --(.whl) file

Decision Tree:

Supervised learning set

"WEKA" \rightarrow DM Tool
New Zealand Waikato University

x_0	x_1 feet ²	x_2	x_3	x_4	20.05.18.
		# bedrooms	# floor	Age	Price
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1535	3	2	30	315
1	1535 (estimate)	2	1	36	178
1.	852				

$n = \# \text{ features}$

$$x^{(1)} \quad x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \rightarrow x_1, x_2, x_3, x_4$$

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

$$x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \\ x_4^{(i)} \end{bmatrix} \in \mathbb{R}^4$$

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

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

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

let, $x_0 = 1$ for all $i = 1$ to m

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

$$= \sum_{j=0}^n \theta_j x_j \quad \rightarrow \theta^T x \\ = \theta^T x \quad (\text{no. of features})$$

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

$n+1$ dimensional vector

$(n+1) \times 1$ matrix x

$$\theta^T x = [\theta_0 \ \theta_1 \ \dots \ \theta_n] \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

* Gradient Descent:

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 \quad (\text{no. of examples})$$

$$= \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$$

when, $n = 1$ # follow the first example

Repeat {

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

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

}

when $n \geq 2$

Repeat {

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

}

Update simultaneously

~~Assignment
Derivation of
Next class~~

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

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

⋮

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

Gradient Descent in practice:

→ Feature Scaling

→ Mean normalization

→ Debug

→ select α

$$x_1 = \text{size} = 100 - 200 \text{ feet}^2$$

$$x_2 = \# \text{bedrooms} = 2 - 5$$

* θ will descent quickly on small range

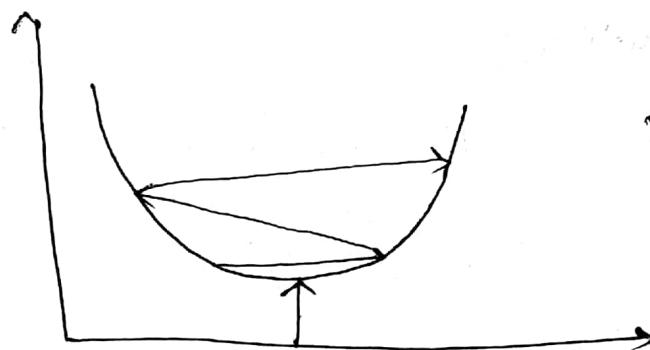
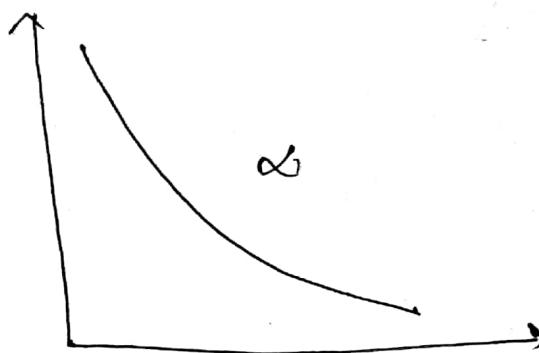
So, set input value in same range (roughly),

$$-1 \leq x^{(i)} \leq 1$$

$$-0.5 \leq (x^{(i)}) \leq 0.5$$

$$[-3 \leq x^{(i)} \leq 3]$$

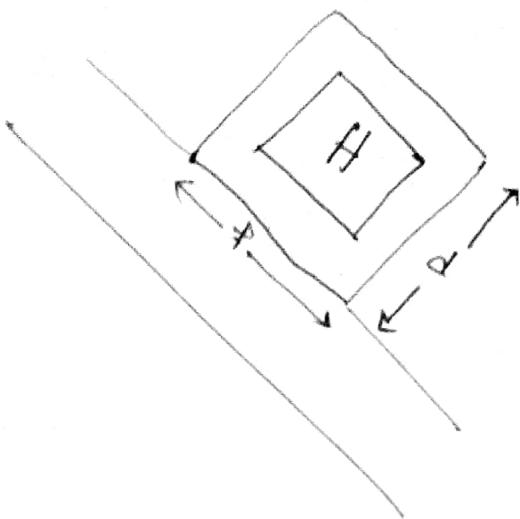
Select α :



$$\alpha = 0.001, 0.003, 0.01, 0.03$$

* for sufficient small α , $J(\theta)$ increase for every iteration.

22.05.18.



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

\uparrow

f

\uparrow

d

\uparrow

x_1

\uparrow

area

\uparrow

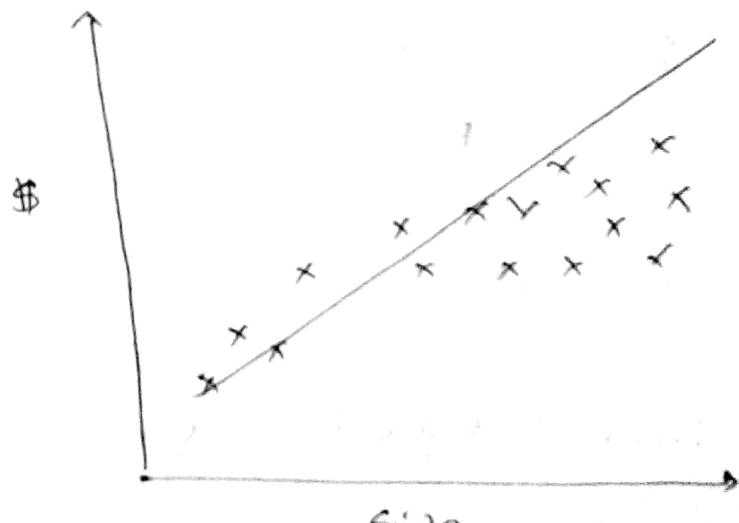
$$= \theta_0 + \theta_1 x$$

*Polynomial Regression

$$x_1 = \text{size}$$

$$x_2 = \text{size}^2$$

$$x_3 = \text{size}^3$$



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

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

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

Feature Scaling:

$$1 - 1000$$

$$x_1 = 1000$$

$$x_2 = 1000000$$

$$x_3 = 10^9$$

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

Normal Equation:

x_0	feet ²	#bedrooms	#floor	Age	Price
1	2104	5	1	45	316
1	1416	3	2	40	420
1	1535	3	2	30	290
1	852	2	1	35	330

* Normal Equation gives other way of minimizing $J(\theta)$. Here we minimize $J(\theta)$ explicitly taking its derivatives w.r.t. the θ_j and setting them to zero. This allows us to find the optimal θ without iteration.

$$\boxed{\theta = (X^T X)^{-1} X^T y}$$

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} = \mathbb{R}^{n+1}$$

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \text{output}$$

X = input matrix

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1535 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 35 \end{bmatrix}$$

$$x^{(1)} = \begin{bmatrix} 1 \\ 2104 \\ 5 \\ 1 \\ 45 \end{bmatrix}$$

$$x^{(2)} = \begin{bmatrix} 1 \\ 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$$X = \begin{bmatrix} x_i^{(1)^T} \\ x_i^{(2)^T} \\ \vdots \\ x_i^{(n)^T} \end{bmatrix}$$

$$X \quad X^T$$

$m \times (n+1) \quad (n+1) \times m \rightarrow \text{dimensions}$

4×5

Feature scaling not needed

Gradient Descent

1. Need to choose α .

2. Need many iteration.

3. $O(kn^2)$

4. Works well if n is

large

Normal Equation

1. No need to choose α .

2. No iteration needed

3. $O(n^3) [(X^T X)^{-1}]$

[if no. of feature is huge,
then complexity would
increase]

4. Slow if n is large

[n = no. of features]

* If $n \geq 1000$ it is better to use gradient descent

$$(X^T X)^{-1}$$

not invertible,

or regularization

1. Redundant features

2. Delete feature
 $n \geq m$

* Assignment → finding weight from height using Python. Before → 5th June.

Sklearn linear regression

udacity

Introduction to machine learning.

27.05.18

Classification :

Email : Spam / Not

Online Tran. : Normal / Not

Tumour : malignant / Benign

Binary classification : $y \in \{0, 1\}$

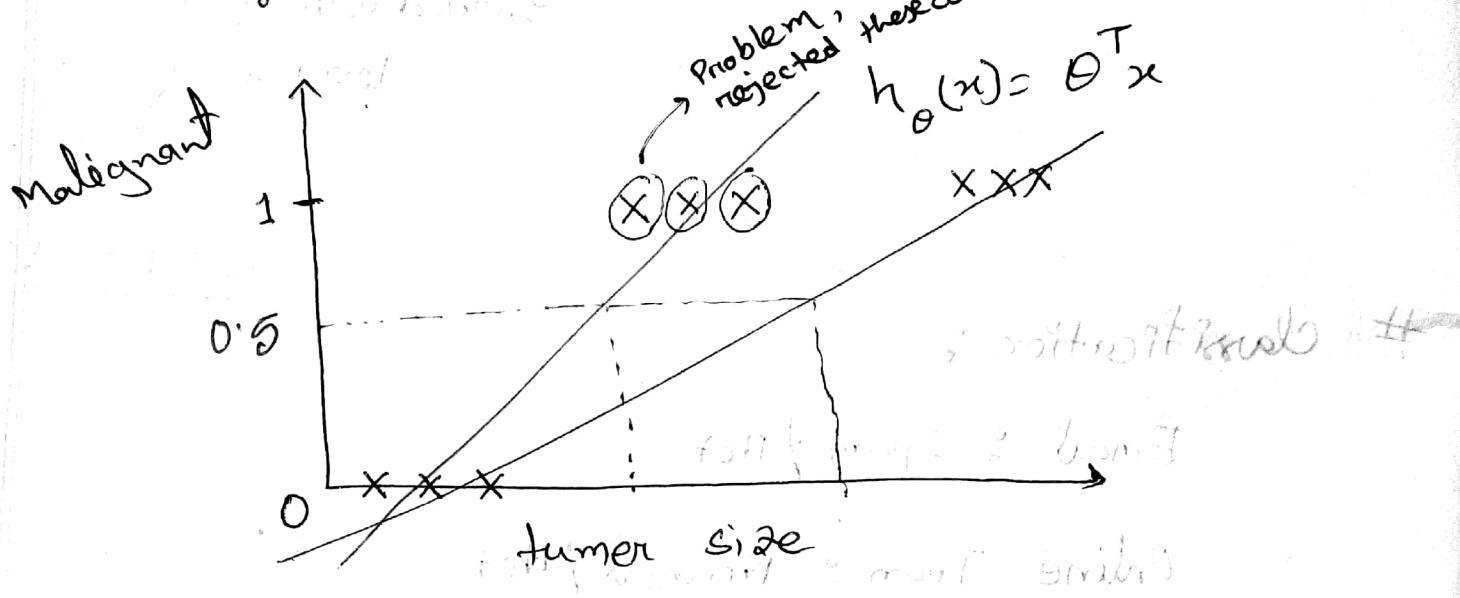
Multiclass Classification : $j \in \{0, 1, 2, 3, \dots\}$

0 : Negative class / Absence of something

1 : Positive class / Presence of something

* Problem with linear Regression for classification

Problem:



$$h_{\theta}(x) \geq 0.5, y = 1$$

$$h_{\theta}(x) < 0.5, y = 0$$

Example of a 2-class classification algorithm

Problem:

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

$h_{\theta}(x)$ using Linear Regression may produce value
 $< 0 \text{ or } > 1$

* Solution: →

Logistic Regression (Logistic function, Sigmoid function)
[Classification algorithm]

$$\frac{1}{1 + e^{-x}}$$

always produces result -

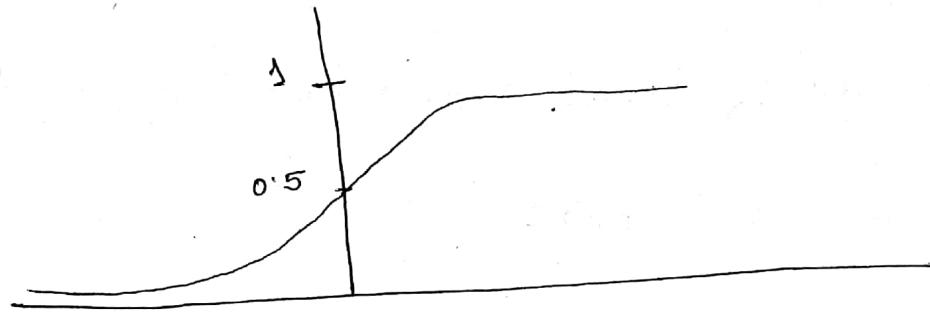
$$0 \leq h_{\theta}(x) \leq 1$$

Hypothesis:

$$h_{\theta}(x) = g(z) ; z = \theta^T x / \theta_0 + \theta_1 x$$

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

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



$$g(z)$$

$$g(0) = 0.5$$

$$1 - e^{-z}$$

$$g(\infty) = 1$$

$$g(-\infty) = 0$$

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

$h_{\theta}(x)$ → Estimate Probability that $y=1$ on input x .

$$h_{\theta}(x) = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{turner-size} \end{bmatrix} = P(y=1 | x, \theta) = h_{\theta}(x) = 0.7$$

$h(x) = P(y=1 | x, \theta)$ Probability that $y=1$,
 $y=0$ or 1 given x parameterized
 by θ

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

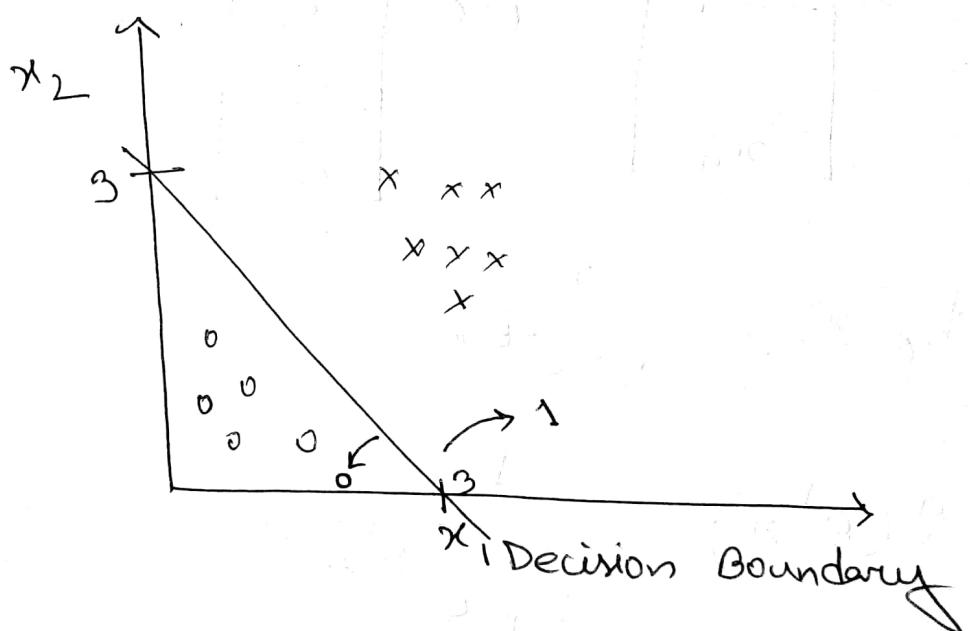
$$P(y=1 | x, \theta) = 1 - P(y=0 | x, \theta)$$

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

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

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

Decision Boundary is the line that separates the area where $y = 0$ and $y = 1$. It is created by hypothesis function.



$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 = \theta^T x$$

$$-3 + 1 x_1 + 1 x_2$$

$$\begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$-3 + x_1 + x_2 \geq 0$$

$$-3 + x_1 + x_2 = 0$$

$$x_1 + x_2 = 3$$

— x —

29.05.18.

Logistic Regression:

Training Examples: $\{(x^{(i)}, y^{(i)})\}, \dots, (x^{(m)}, y^{(m)})\}$ # Example: m # Features: n

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \theta \in \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

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

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

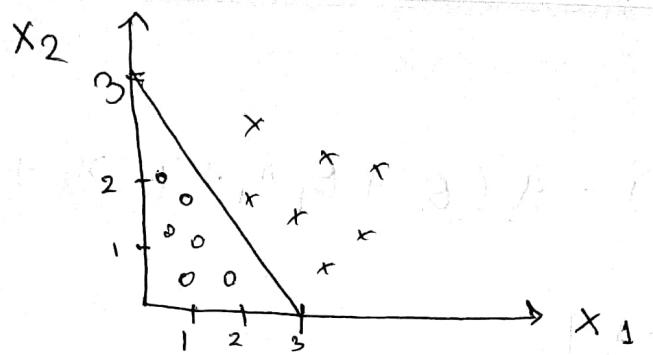
$$Y = 1, \quad g(2) \geq 0.5$$

$$\theta^T x = 0 \nearrow$$

$$g \neq \theta^T x \geq 0, \quad Y = 1$$

$$\text{else } Y = 0$$

2



$$\theta^T x \rightarrow h_{\theta}(x)$$

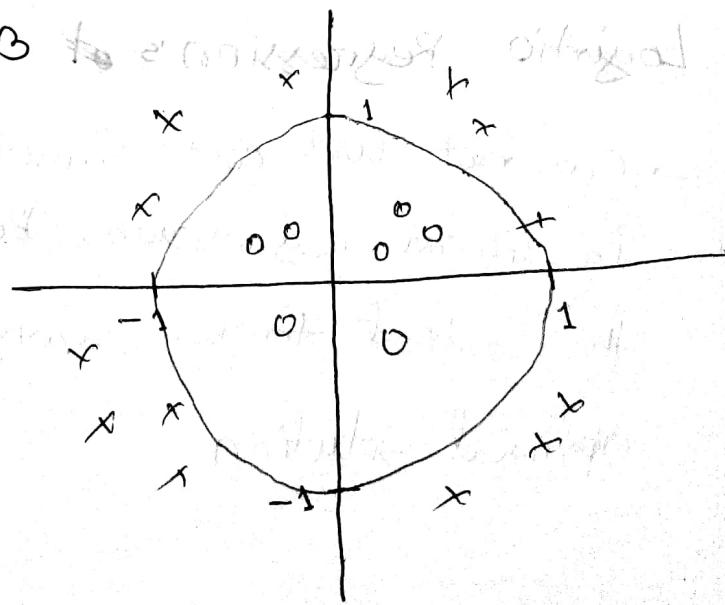
$$\theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$\begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} = -3 + x_1 + x_2 \quad \left| \begin{array}{l} Y=1 \text{ when } \theta^T x \geq 0 \\ Y=0 \text{ when } \theta^T x < 0 \end{array} \right.$$

$-3 + x_1 + x_2 \geq 0$ is another notation

$x_1 + x_2 \geq 3$, the graph for this

$-3 + x_1 + x_2 = 3$ to consider straight



3

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

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

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} = -1 + x_1^2 + x_2^2$$

$$-1 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 = 1 \quad [\text{equ. for circle}]$$

* Decision boundary is a property of hypothesis, not of training set.

Logistic Regression's cost function:

→ Can not use cost function that we've used in linear regression. Because, it will cause the output to be wavy, causing many local optimal solution.

4

→ it will not be ~~not~~ convex function.

* Cost function:

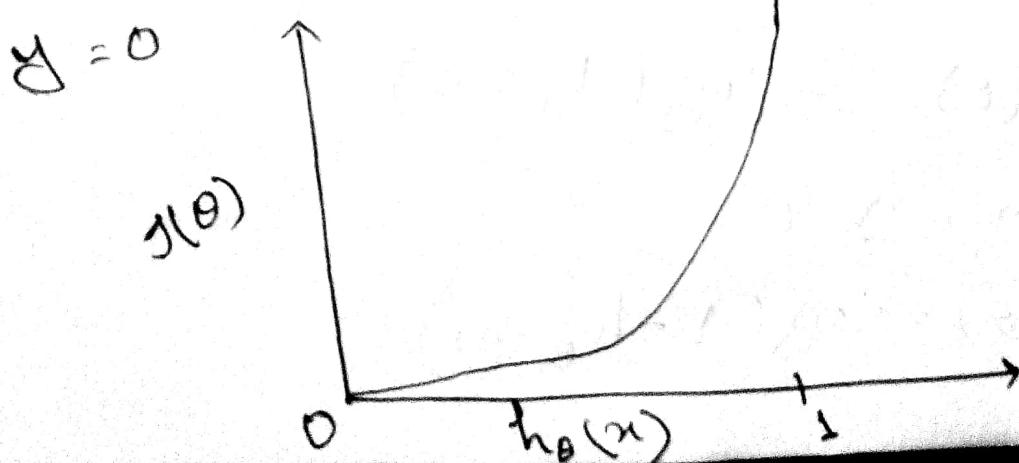
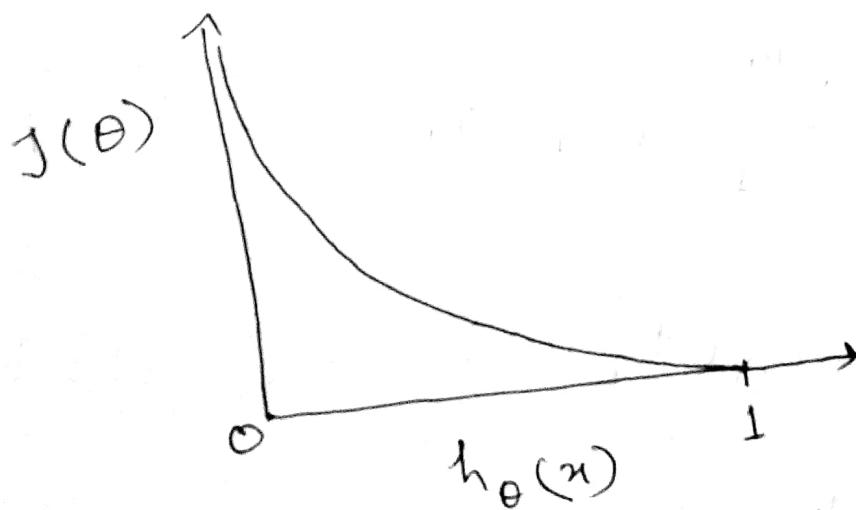
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}\left(g\left(\sum_{i=0}^n \theta^T x^{(i)}\right), y^{(i)}\right)$$

$y = 1$ or 0

when

$$\hat{y} = 1$$



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

$$J(\theta) = 0, \quad y = h_{\theta}(x)$$

$$J(\theta) = \infty, \quad y = 1, \quad h_{\theta}(x) = 0$$

$$J(\theta) = \infty, \quad y = 0, \quad h_{\theta}(x) = 1$$

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

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

$$\text{when, } y = 1$$

$$J(\theta) = -\log(h_{\theta}(x))$$

$$\text{when, } y = 0$$

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

6

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

$$J(\theta) = \frac{1}{m} \left(-\hat{y}^T \log(h) - (1-\hat{y})^T \log(1-h) \right)$$

Gradient Descent:

Repeat {

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

// update
simultaneously

}

Repeat {

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

}

Vector form: $\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \hat{y})$

• OFGIS

• L-OFGIS

7

03.06.18.

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



$$\hat{y}(z)$$

0.5

$$Y = 1$$

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

$$\begin{cases} 1 & -\log(h_{\theta}(x)) \\ 0 & -\log(1 - h_{\theta}(x)) \end{cases}$$

↓
 $g(z)$

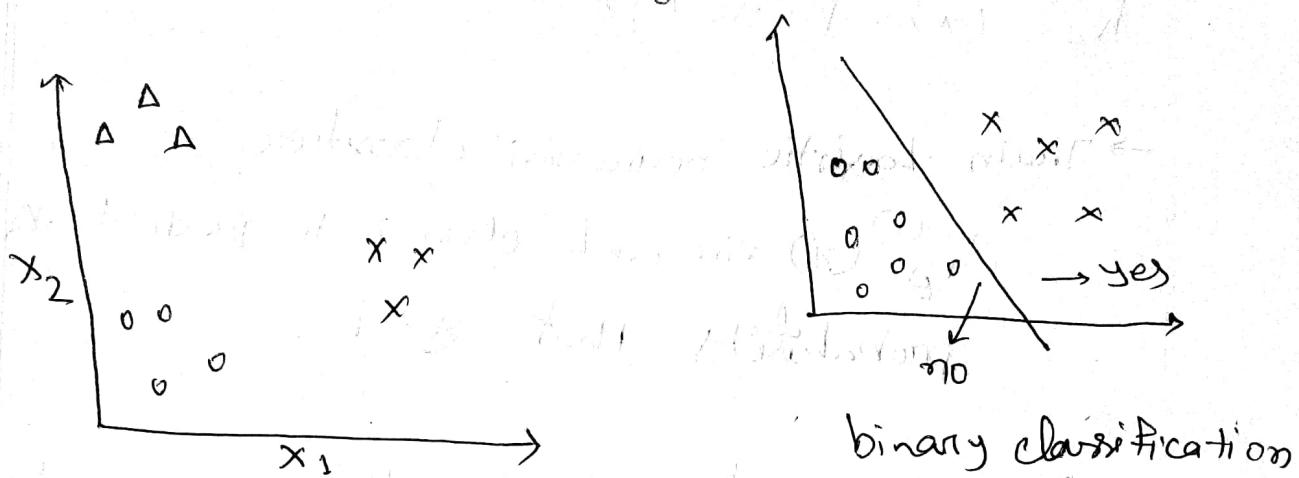
$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m (y_i \log(h_{\theta}(x)) + (1-y_i) \log(1 - h_{\theta}(x)))$$

$Y = 1, h_{\theta}(x) \rightarrow 0 ; \text{cost} = \text{infinity}$

$Y = h_{\theta}(x) ; \text{cost} = \text{zero}$

8

Multiclass classification



Multiclass classification

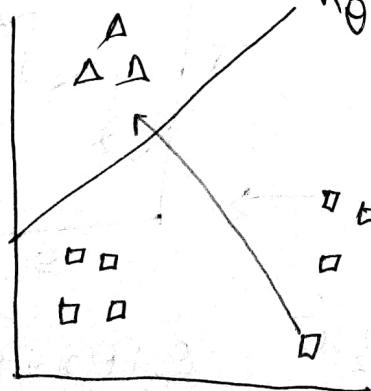
* One vs all (one vs rest)

class 1: Δ

class 2: O

class 3: X

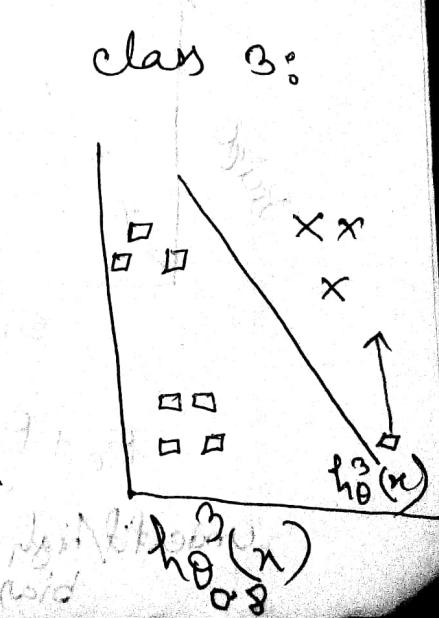
Class 1: $h_0^1(x)$



Class 2: $h_0^2(x)$



Class 3:



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

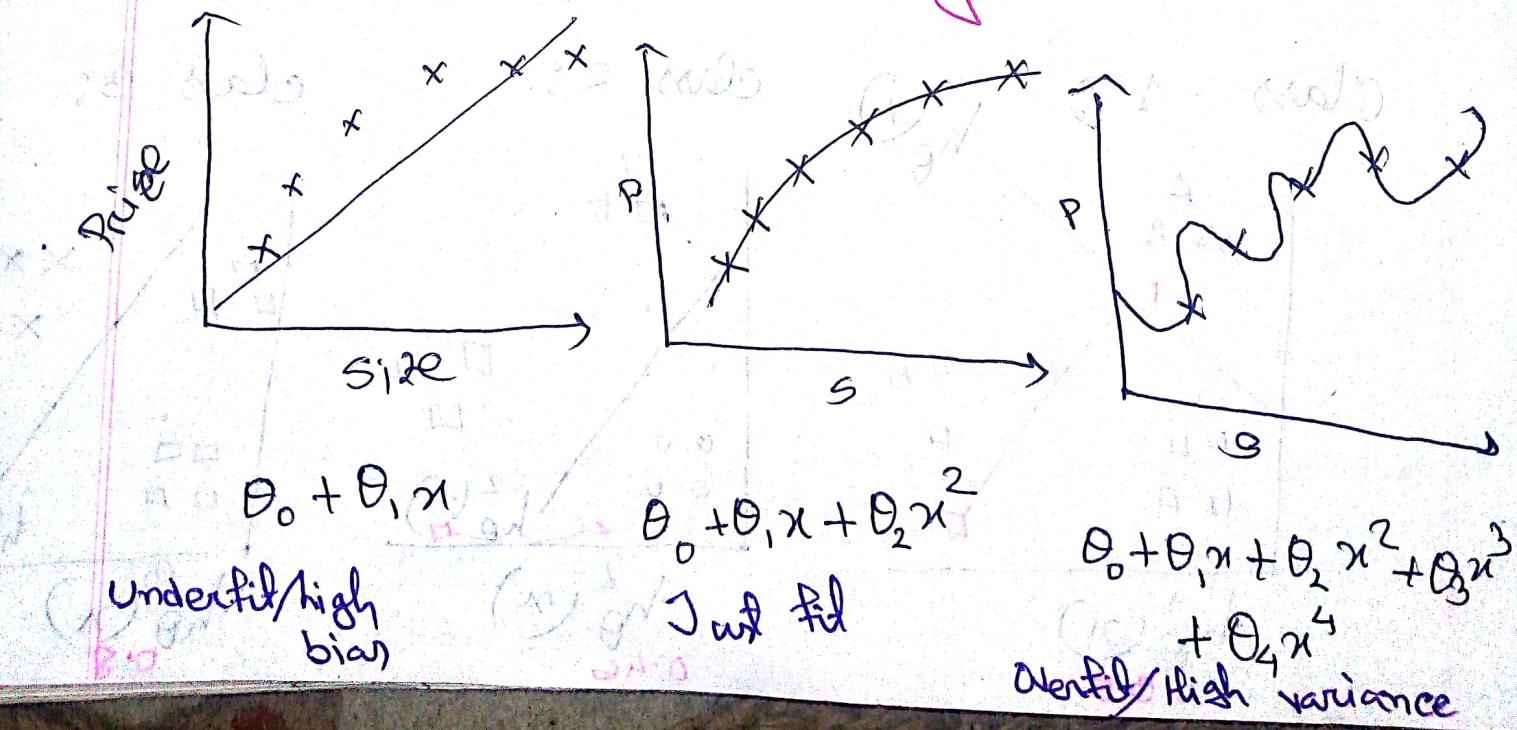
→ Train Logistic Regression classifier.

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

→ One new input x , to make a prediction
pick the class i^{th} maximize

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

Overfitting Vs Underfitting



Underfit/high bias

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

Just fit

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

Overfit/high variance

* Underfit:

Structure not captured by model.

$h_{\theta}(x)$ maps poorly.

* Overfit:

Fits training data set very well ($J(\theta) \approx 0$),
but fails to generalize to new input

Addressing overfitting problems

1. Reduce the # of features

→ manually select which feature to keep.

→ use a model selection algorithm

2. Regularization

→ keep all features, but reduce the magnitude of θ_s .

→ works well when we have a lot of slightly useful features.

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

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

↓ ↓
tends to $\theta \rightarrow 0$ 0

$$\begin{matrix} x_0 & x_1 & x_2 & x_3 & \dots & x_n \\ \theta_0 & \theta_1 & \theta_2 & \theta_3 & \dots & \theta_n \end{matrix} \left. \begin{array}{l} \text{don't know} \\ \text{which one} \\ \text{to reduce.} \end{array} \right\}$$

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

; λ = regularization constant.

— x —

R

05.06.18.

Regularization: (for linear Regression)

 \Rightarrow small value for θ 's \rightarrow simple hypothesis \rightarrow less prone to overfitting

if λ high,
 θ will be
 small.

Features: $x_0, x_1, x_2, \dots, x_n$ P: $\theta_0, \theta_1, \theta_2, \dots, \theta_n$

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

Gradient descent:

Repeat {

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

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

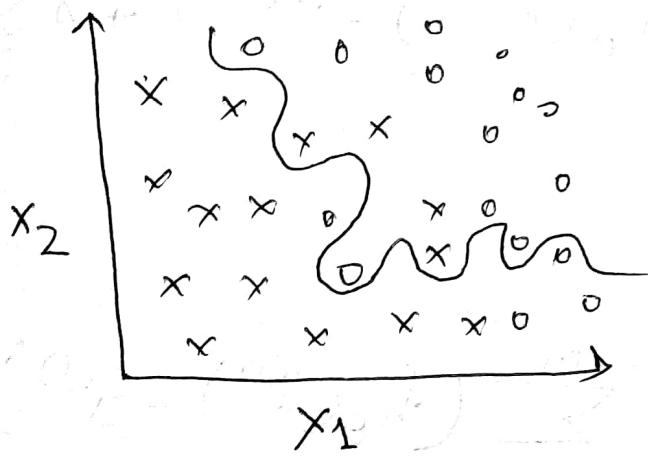
}

13

$$\theta = (X^T X + \lambda L)^{-1} X^T Y ; L = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{n+1}$$

 $\rightarrow x$

Regularization for Logistic Regression



$$h_{\theta}(x) = g(\underline{\theta})$$

$$= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$+ \theta_5 x_1 x_2 + \theta_6 x_1^2 x_2 + \dots)$$

$$\begin{aligned}
 * \text{ Cost } \% \quad J(\theta) = & - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + \right. \\
 & \left. (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)})) \right] \\
 & + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2
 \end{aligned}$$

Gradient descent:

repeat {

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

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

}

$\rightarrow x \rightarrow$

Function Linear-Reg (theta_0, theta_1, x)

hy^o

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

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

GD^o

repeat {

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

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

}

Fun Linear (theta_0, theta_1 , x)

return $\theta_0 + \theta_1 * x$

costfun (x)

for i = 1 to n

$$\text{cost} = \text{cost} + (\text{Linear}(\theta_0, \theta_1, x^{(i)}) - y^{(i)})^2$$

return cost * $\frac{1}{2m}$

$x[80] = \{ \dots \}$ // global
 $y[80] = \{ \dots \}$

GDfun (θ_0, θ_1)

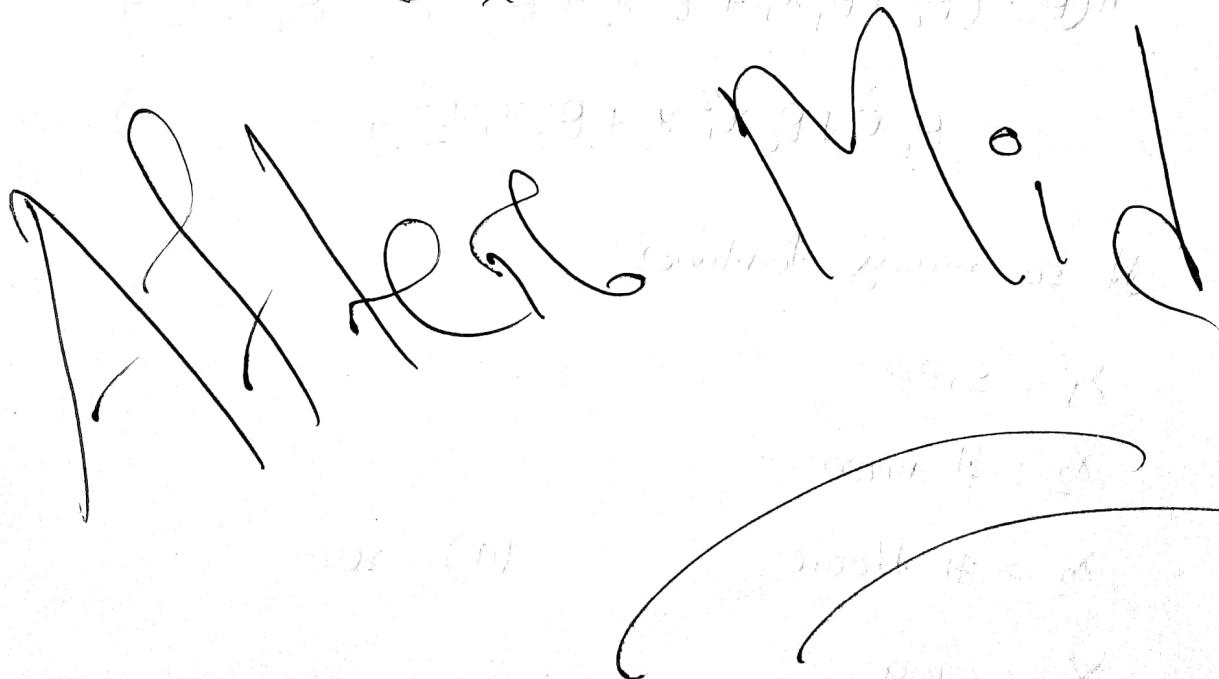
for $i = 1$ to n // repeat 1000 times

$$\theta_0 = \theta_0 - (x^{(i)} - y^{(i)})$$

$$\theta_1 = \theta_1 - (x^{(i)} - y^{(i)}) x^{(i)}$$

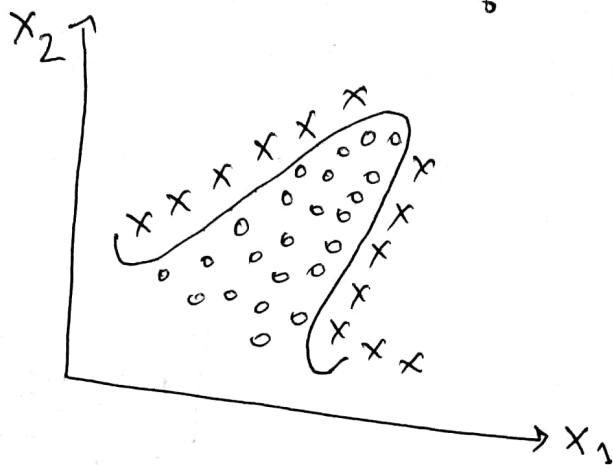
Kaggle

registration.



Neural Networks

* classification problems



$$\vartheta(z)$$

↓

$$h(\theta) = (\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1 x_2 + \theta_5 x_1^2 + \theta_6 x_1^2 x_2 + \theta_7 x_1 x_2^2 + \dots)$$

if so many features

$$x_1 = \text{size}$$

$$x_2 = \# \text{ room}$$

$$x_3 = \# \text{ floor} \quad (n) = 100$$

$$x_4 = \text{area}$$

$$\dots$$

$$x_{100} = \dots$$

$$x_{100} =$$

2nd order equation is needed,

$$x_1^2 + x_1 x_2 + x_1 x_3 + \dots + x_1 x_{100} + x_2^2 + x_2 x_3 + x_2 x_4 + \dots + x_2 x_{100} + \dots$$

Complexity $O(n^2)$

Gray image \rightarrow 8 byte

Color image \rightarrow 32 bytes

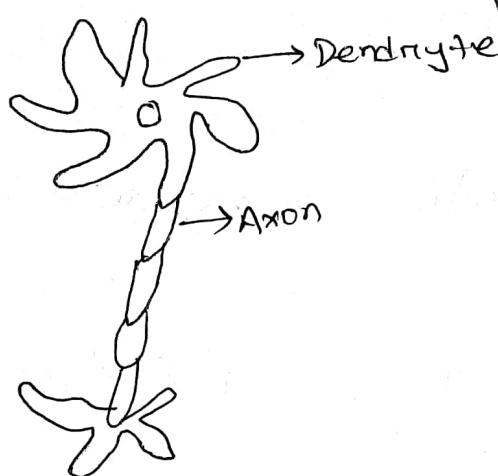
\rightarrow -

Least weight coefficients

Artificial Neural Network

26.06.18.

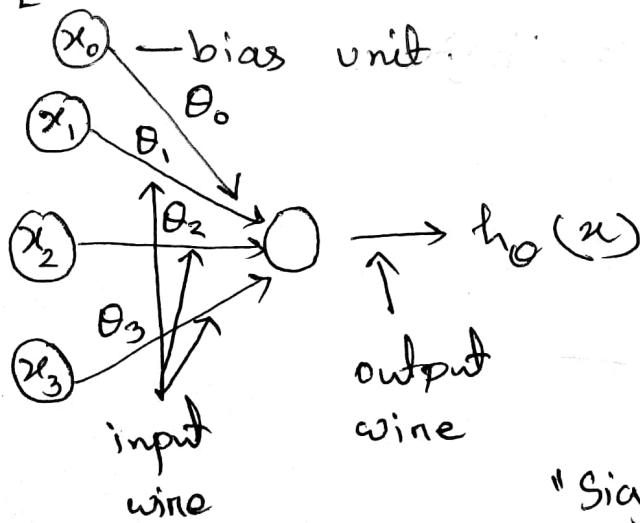
(for classification problem)



$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

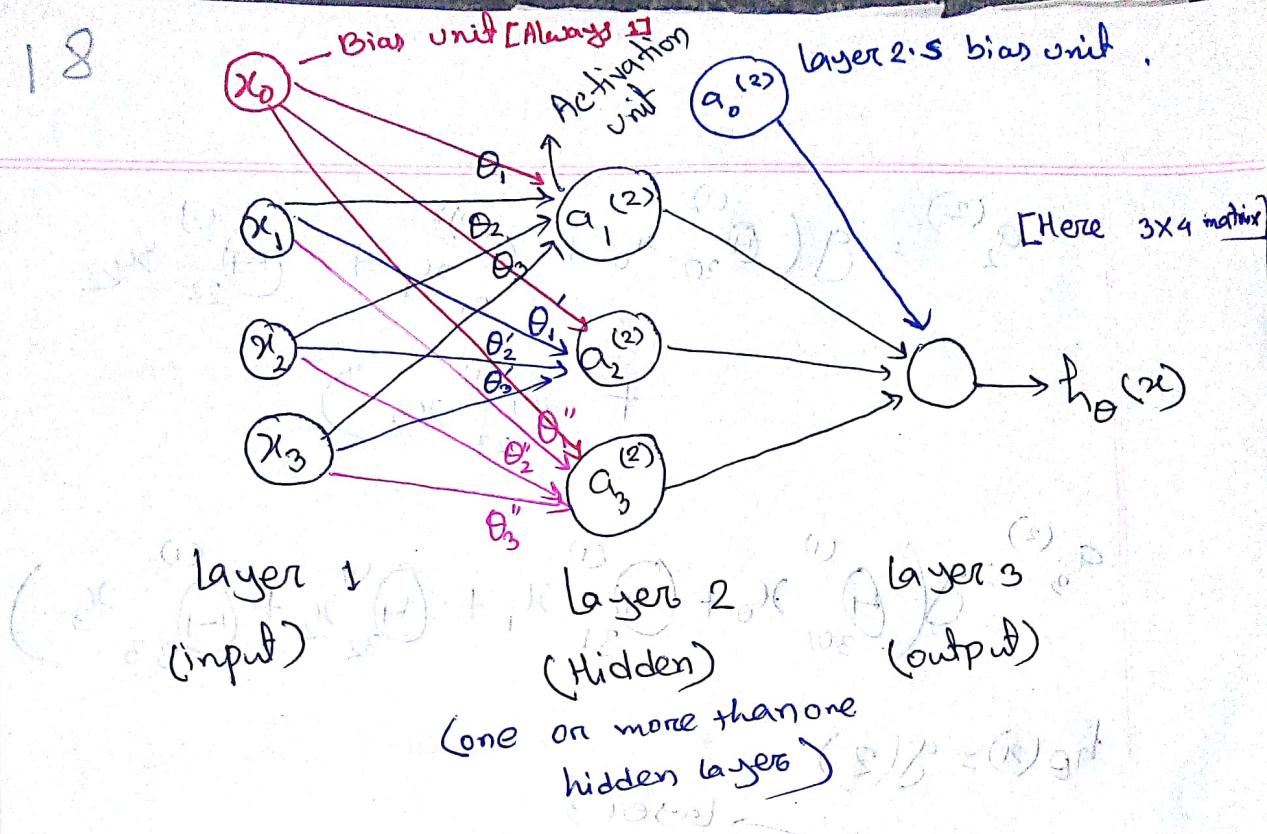
Parameter/weight.



"Sigmoid Activation function"

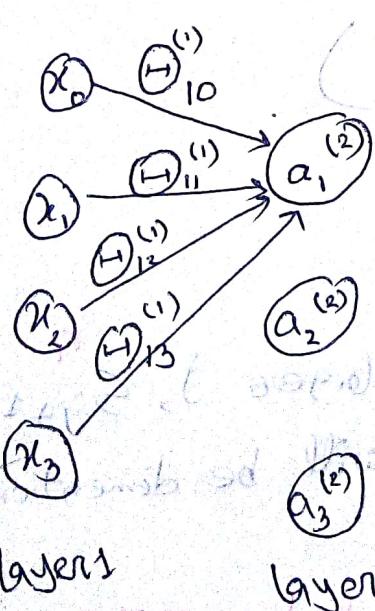
$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{-\theta^T x}} \quad (\text{use linear order})$$

18



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

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



$$\begin{aligned}
 a_1^{(2)} &= g(z) \\
 &= g(\theta^T x) \\
 &= g(\theta_{10} x_0 + \theta_{11} x_1 + \theta_{12} x_2 + \theta_{13} x_3)
 \end{aligned}$$

$$(1 + e^{-z}) \times 1 + e^z$$

per XMR (soft, tanh and sigmoid) on $f(x)(1+m)x^n$

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

l.e.o

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

$$h_\theta(x) = g(z)$$

→ layer 3

$$a_1^{(3)} \rightarrow \text{Unit}_1$$

$$= g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)}$$

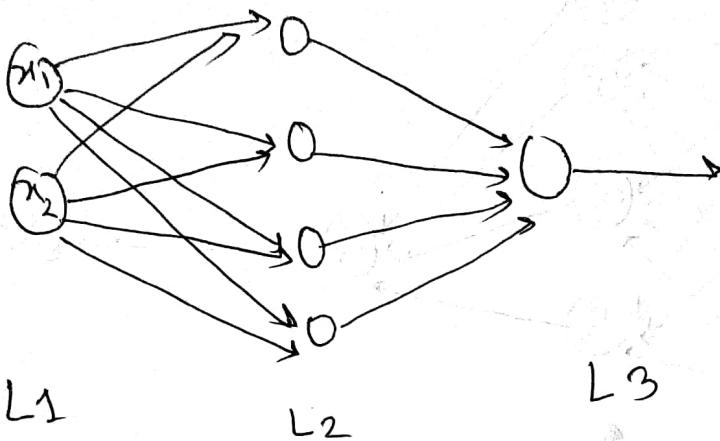
$$+ \theta_{13}^{(2)} a_3^{(2)})$$

$$\Theta^1 \in \mathbb{R}^{3 \times 4}$$

* If a network has s_j^m units in layer j , s_{j+1}^n units in layer $j+1$, then $\Theta^{(j)}$ will be dimension of $s_{j+1}^n \times (s_j^m + 1)$
 $n \times (m+1)$; if we calculate bias unit, else $n \times m$.

20

Ques.



if we use bias unit, Dimension would be
 $(4 \times (2+1))$

$$\text{Ans} = 4 \times (2+1) \times 3 = 4 \times 3$$

$$(\text{Ans} \text{ else } 4 \times 2 \times 3 + 5 \times 3) \times B = 4 \times 3$$

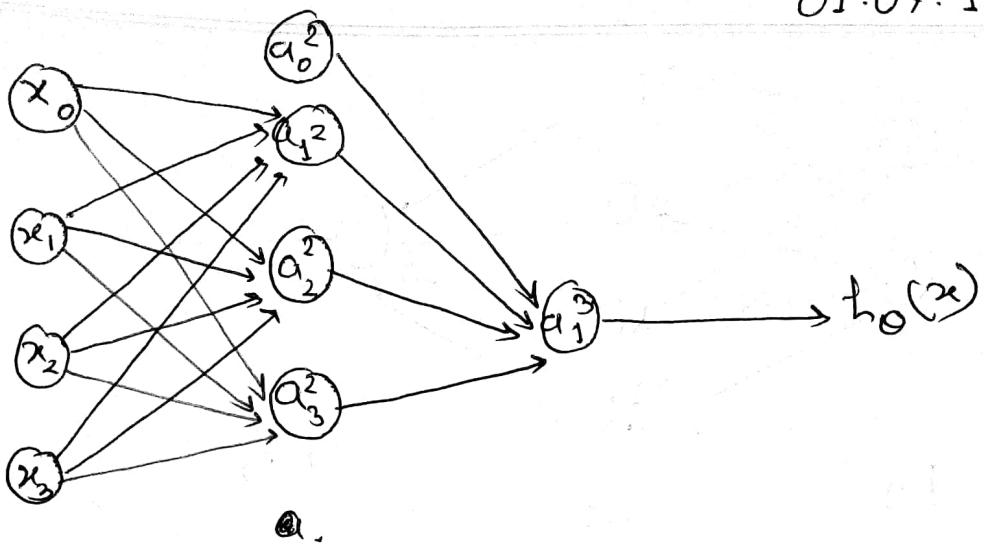
(Ans)

Ans without bias unit

$$(4 \times 3) \times B = 4 \times 3$$

$$(4 \times 2) \times B = 4 \times 2$$

$$(4 \times 1) \times B = 4 \times 1$$



$$a_1^{(2)} = g(\theta_{1,0}^1 x_0 + \theta_{1,1}^1 x_1 + \theta_{1,2}^1 x_2 + \theta_{1,3}^1 x_3)$$

$$a_2^{(2)} = g(\theta_{2,0}^1 x_0 + \theta_{2,1}^1 x_1 + \theta_{2,2}^1 x_2 + \theta_{2,3}^1 x_3)$$

$$a_3^{(2)} = g(\theta_{3,0}^1 x_0 + \theta_{3,1}^1 x_1 + \theta_{3,2}^1 x_2 + \theta_{3,3}^1 x_3)$$

$$a_1^{(3)} = g(\theta_{1,0}^{(2)} a_0^{(2)} + \theta_{1,1}^{(2)} a_1^{(2)} + \theta_{1,2}^{(2)} a_2^{(2)} + \theta_{1,3}^{(2)} a_3^{(2)})$$

using logistic function $g(z)$

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

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

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

22

$$\underline{z}^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$\underline{z}^{(2)} = \frac{\underline{x}}{\mathbb{R}^3} \circled{H}^1$$

3x4

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

$$\mathbb{R}^3 \rightarrow \mathbb{R}^3$$

For layer $j = 2$ and k unit,

$$z_k^{(2)} = \theta_{k0}^{(2)} x_0 + \circled{H}_{k1}^{(2)} x_1 + \dots$$

$$\underline{z}^{(2)} = \circled{H}^{(2)} \underline{x}^{(2)}$$

$$\mathbb{R} \times \mathbb{R}^{3 \times 4}$$

$$\underline{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix}$$

$$\underline{z}^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix} = \underline{z}^{(j-1)} \circled{H}^{(j)}$$

$$\Rightarrow a^j = g(z^{(j)})$$

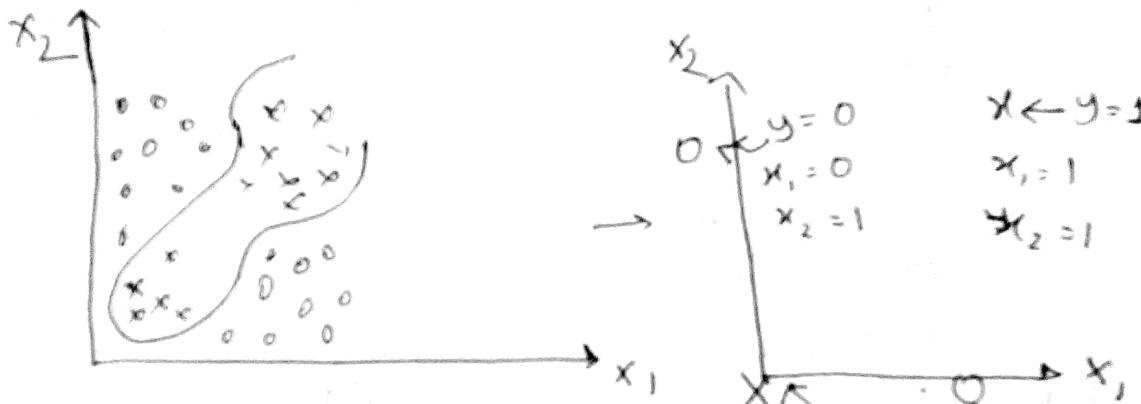
23

$$a_1^3 = g(2_1^3)$$

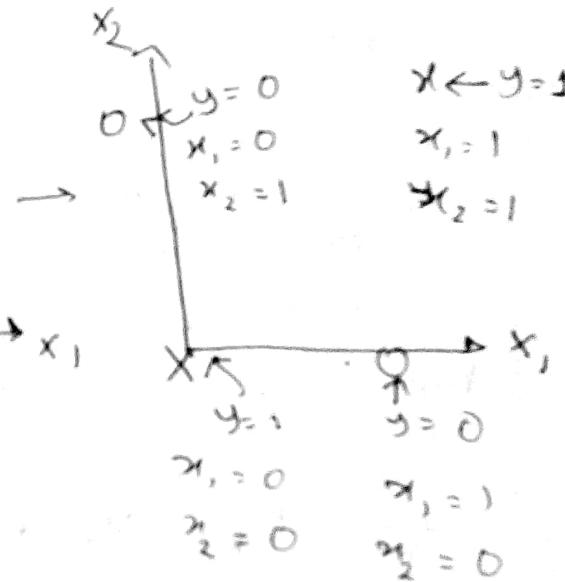
$$\rightarrow \frac{1}{1+e^{-2_1^3}}$$

(logistic Regression)

Q: Why neural network is imp to design non-linear hyp.?



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



x_1 on x_2

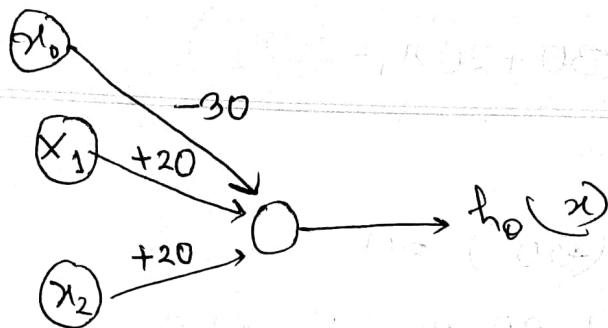
NOT (x_1 XOR x_2)

x_1 XNOR x_2

AND, OR, NOT

x_1	x_2	$x_1 \oplus x_2$	$g(x_1 \oplus x_2)$
0	0	0	1
0	1	1	0
1	0	1	0
1	1	0	1

24



$$h_0(x) = g(x^2)$$

$$= g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2)$$

$$= g(-30 + 20x_1 + 20x_2)$$

$$\begin{matrix} x_1 \\ 0 \end{matrix}$$

$$\begin{matrix} x_2 \\ 0 \end{matrix}$$

$$\rightarrow g(-30 + 20 \cdot 0 + 20 \cdot 0) = g(-30) = 0$$

$$\begin{matrix} x_1 \\ 0 \end{matrix}$$

$$\rightarrow g(-30 + 20) = g(-10) = 0$$

$$\begin{matrix} x_1 \\ 1 \end{matrix}$$

$$\rightarrow g(-30 + 20) = g(-10) = 0$$

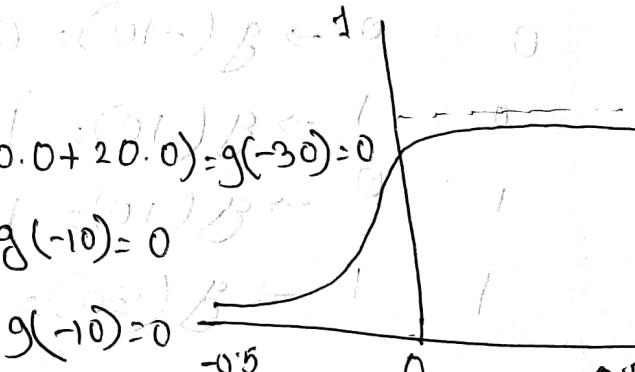
$$\begin{matrix} x_1 \\ 1 \end{matrix}$$

$$\rightarrow g(-30 + 20 + 20) = g(10) = 1$$

x_1 AND x_2

$$g(2) = 1$$

$$z \geq 0$$



$$g(-30 + 20x_1 - 20x_2)$$

x_1 x_2

$$0 \quad 0 \rightarrow g(30) = 0$$

$$0 \quad 1 \rightarrow g(-30 - 20) = g(-50) = 0$$

$$1 \quad 0 \rightarrow g(-30 + 20) = g(-10) = 0$$

$$1 \quad 1 \rightarrow g(-30 + 20 - 20) = (-30) = 0$$

$$g(-10 + 20x_1 + 20x_2)$$

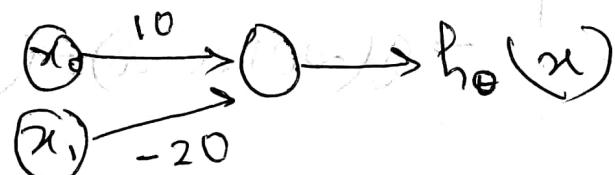
x_1 x_2

$$0 \quad 0 \rightarrow g(-10) = 0$$

$$0 \quad 1 \rightarrow g(10) = 1$$

$$1 \quad 0 \rightarrow g(10) = 1$$

$$1 \quad 1 \rightarrow g(30) = 1$$



x_1

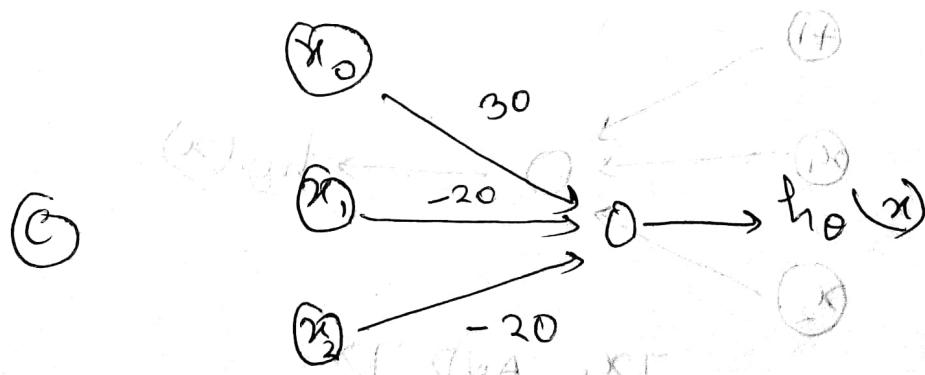
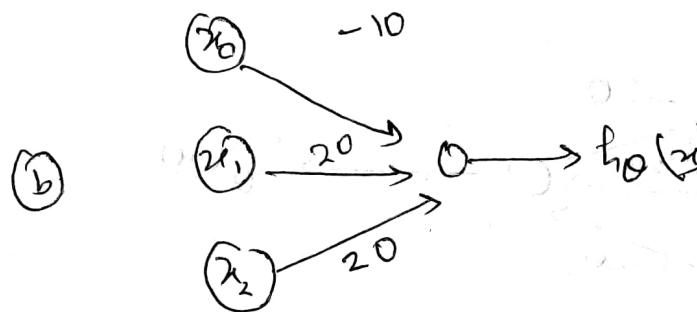
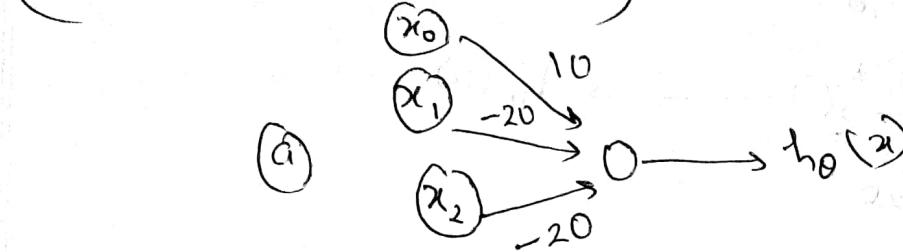
$$0 \rightarrow g(0 - 20) = g(0) \quad g(10) = 1$$

$$1 \rightarrow g(10 - 20) = 0$$

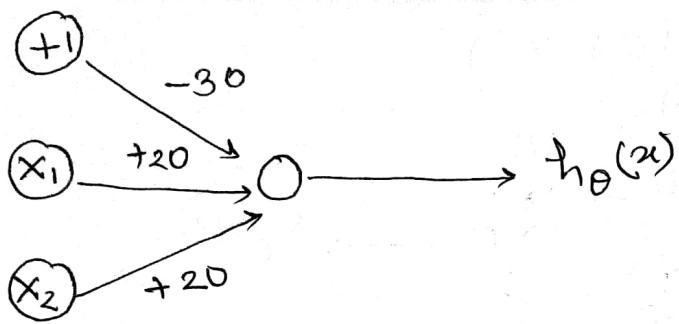
Assignment

26

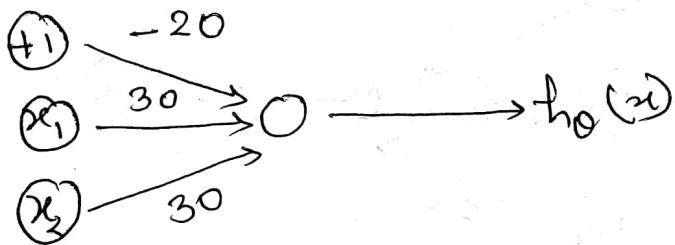
(NOT x_1 AND NOT x_2)



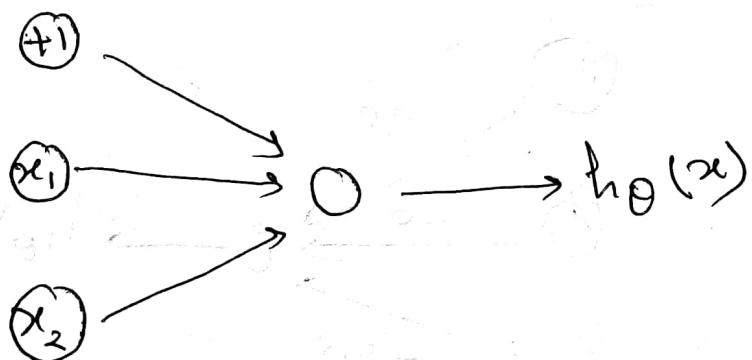
03.07.18. 27



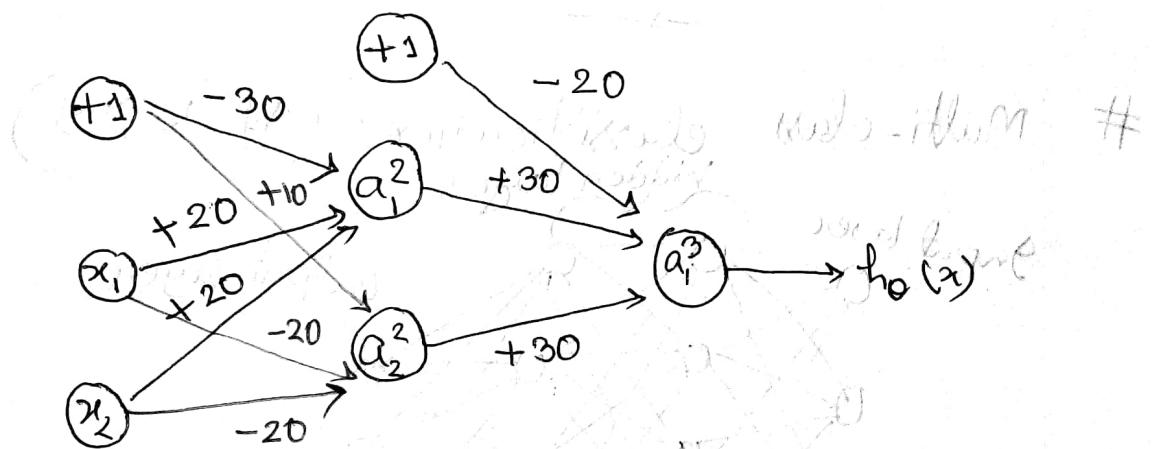
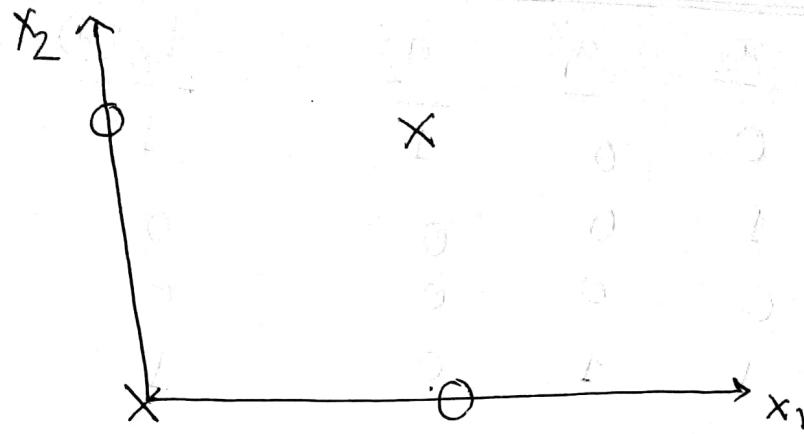
AND



OR



$\neg x_1 \text{ AND } \neg x_2$



(Forward propagation)

$$a_1^2 = 1 \text{ iff } x_1 = x_2 = 1$$

otherwise 0

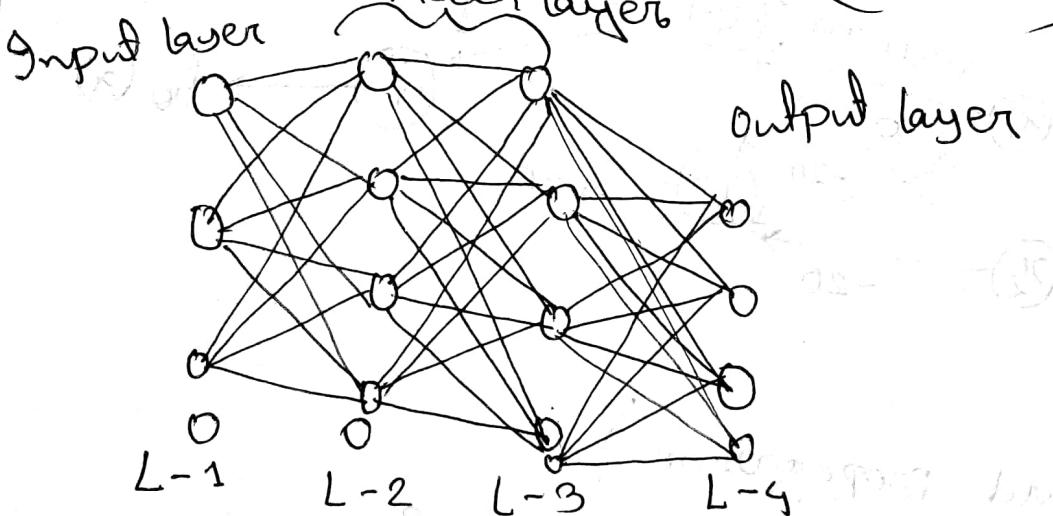
$$a_2^2 = 1 \text{ iff } x_1 = x_2 = 0$$

otherwise = 0

x_1	x_2	a_1^2	a_2^2	$\frac{f_0(x)}{1}$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

 $\rightarrow x$

Multi-class classification: ($f \geq k \geq 3$)



$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

car

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

truck

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Human

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

motor cycle.

30

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(n)}, y^{(n)})\}$$

$L = \# \text{ layers in network}$

$s_l = \# \text{ unit in layer } l$

$$s_1 = 4 \quad s_2 = 5 \quad s_3 = 5 \quad s_4 = 4 = s_L$$

$K = K$ class classification.

$$K = h_\theta(x) \rightarrow \text{vector}$$

$$\rightarrow \mathbb{R}^K$$

Binary classification

$$s_i = ? \quad 1$$

$$h_\theta(x) = 1$$

* If two class, then

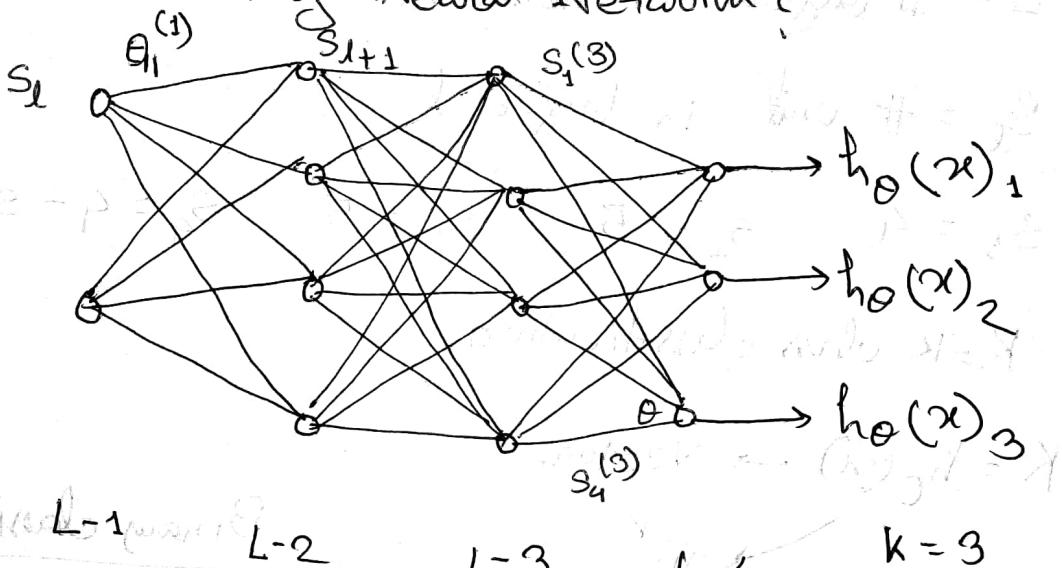
Binary classification.

08.07.18.

Today's Topic:

→ How to find out Cost function and Gradient Descent using Neural Network?

Descent using Neural Network?



* Cost function:

Logistic Regression:

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

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

32

$$h_{\theta}(x) \in \mathbb{R}^k \quad (h_{\theta}(x))_i \text{ is the output unit}$$

* Neural Network:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \left[\sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{0j})^{(l)} \right]^2$$

The double sum adds up the logistic regression costs calculated for each cell in output layer.

The triple adds up the square of all the individual θ s in the entire network, except all $\theta_{0j} \rightarrow$ because i starts with 1, so θ_{0j} will be eliminated.

ith feature

I j^{th} layer to m^{th} layer = (θ) is initiated

j^{th} unit. (θ) the units are not

* Gradient Descent:

min

$$J(\theta) \mid m=1$$

→ example.

Forward propagation,

$a^{(1)}$

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

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

$$a^{(2)} = g(z^{(2)}) \text{ add } (\eta^{(2)})$$

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

$$a^{(3)} = g(z^{(3)}) \text{ along with } \text{ backprop}$$

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

$$a^{(4)} = g(z^{(4)}) = h_{\theta}(x)$$

* Gradient Computation: Back propagation algorithm

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

For each output unit ($l=4$)

34

$$\emptyset \quad \delta_j^{(4)} = a_j^{(4)} - y_j \quad \| \quad \delta^{(4)} = a^{(4)} - y$$

last layer calculates the distance

$$\delta^{(3)} \rightarrow \delta_1^{(3)} \quad \delta_2^{(3)} \quad \delta_3^{(3)} \quad \delta_4^{(3)}$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} * g'(z^{(3)})$$

$g'(z^{(3)})$ is the derivative of activation

function $g(z^{(3)})$

* = element wise multiplication

$$g'(z^{(3)}) = a^{(3)} * (1 - a^{(3)})$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} * a^{(3)} * (1 - a^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} * a^{(2)} * (1 - a^{(2)})$$

$\delta^{(1)}$ → no need to calculate, because input has

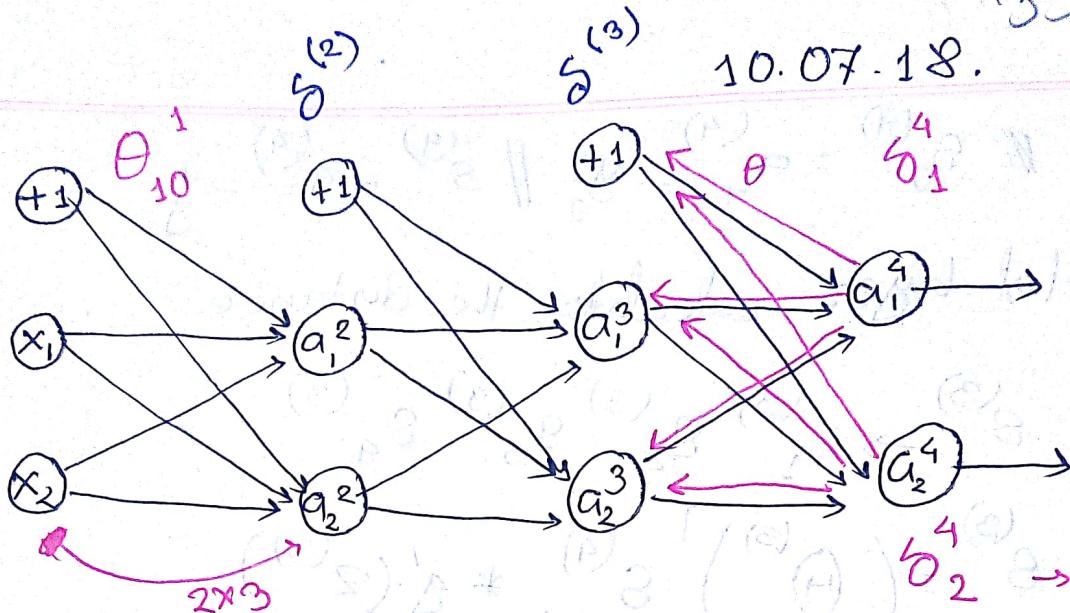
no error

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

[ignoring x_i , if $\lambda=0$].

→ →

10.07.18.

Cost Function $J(\theta)$

$$\frac{\partial J(\theta)}{\partial \theta_j^{(l)}}$$

layer l to layer $l+1$

$\theta_j^{(l)}$ is the weight between unit j in layer l and unit i in layer $l+1$.
 $s_j^{(l)}$ is the error of unit j in layer l .

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

Back Propagation

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

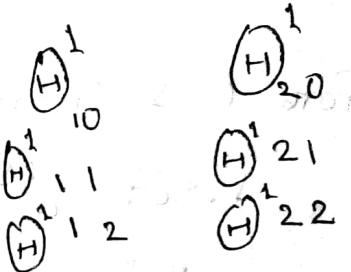
BFSGI

↓ takes

$$J(\theta)$$

& partial derivative
as input

For layer -1



$$\delta^{(4)} = a_j^{(4)} - y_j = a^{(4)} - y$$

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

$$= (\theta^{(3)})^T \delta^{(4)} * a^{(3)} * (1 - a^{(3)})$$

(Back propagation)

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

Back propagation algorithm:

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

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

↳ use to compute $\frac{\partial}{\partial \theta_{ij}} J(\theta)$

- For $t = 1$ to m

1. Set $a^{(1)} := x^{(t)}$

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

3. Using $y^{(t)}$, compute $s^{(L)} = a^{(L)} - y^{(t)}$

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

formula,

$$s^{(l)} = ((\Theta^{(l)})^T \cdot s^{(l+1)}). * \frac{a^{(l)} * (1-a^{(l)})}{g'(z^{(l)})}$$

$$5. \Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} s_i^{(l+1)}$$

Vector form

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

2020-18

Hence we update our new Δ matrix (out of for loop)

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

(Cost function \rightarrow feature scaling to avoid overfitting problem)

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

* The capital-delta matrix for next step
 D is used as an "accumulator" to add up our values as we go along and eventually compute our partial derivative, this we get,

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

To reduce cost function, $J(\theta)$:

→ Gradient Descent

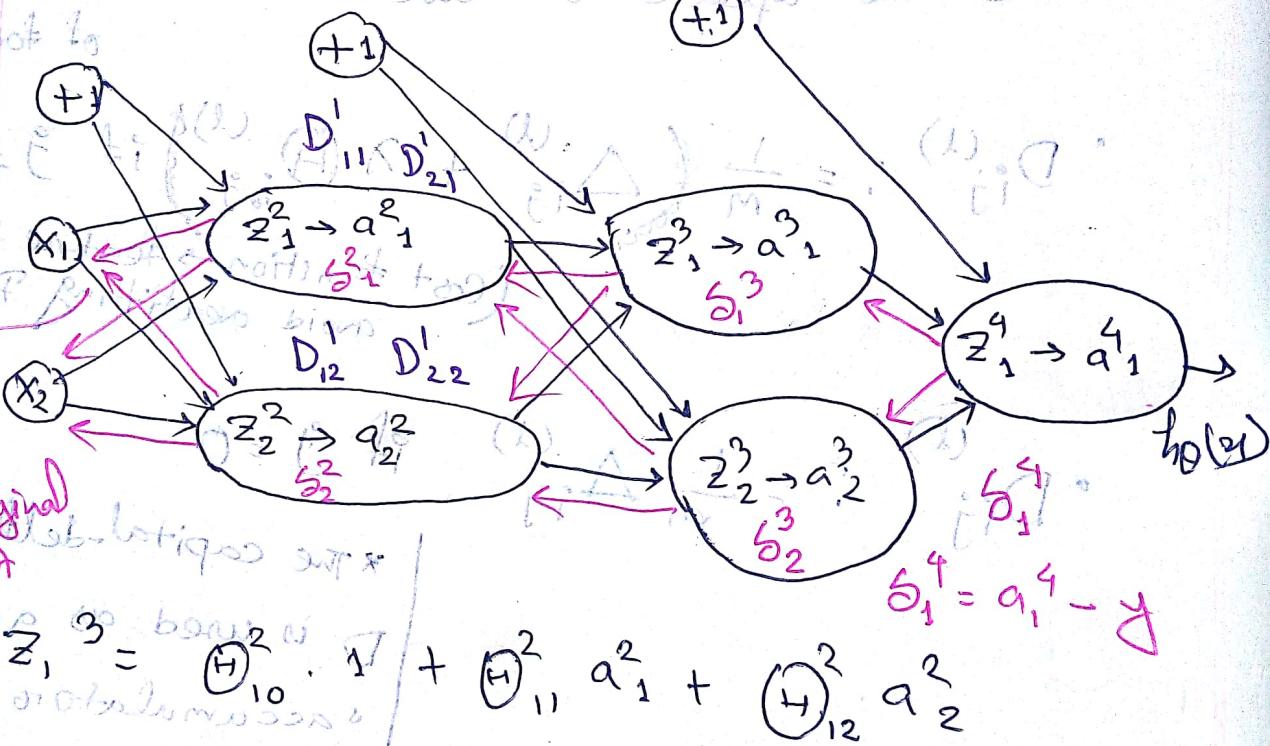
→ Conjugate GD

→ BFGS

→ L-BFGS

(b) \times iteration Δ with two line steps see first

(goal not to)



$$z_1 = H_{10}^{-1} \cdot \nabla + H_{11}^2 a_1^2 + H_{12}^2 a_2^2$$

using

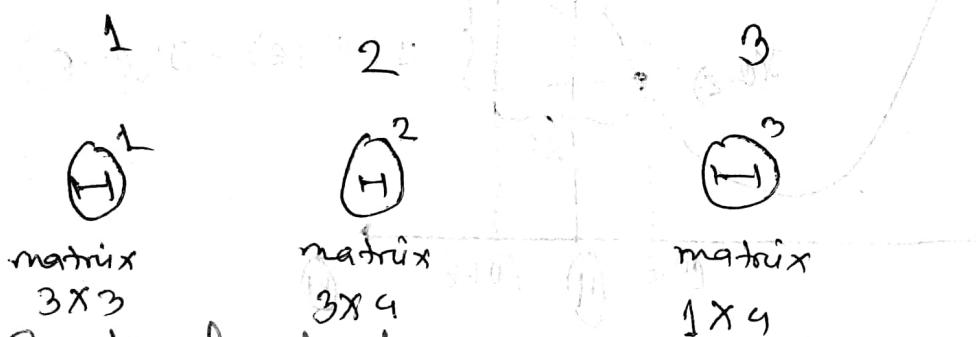
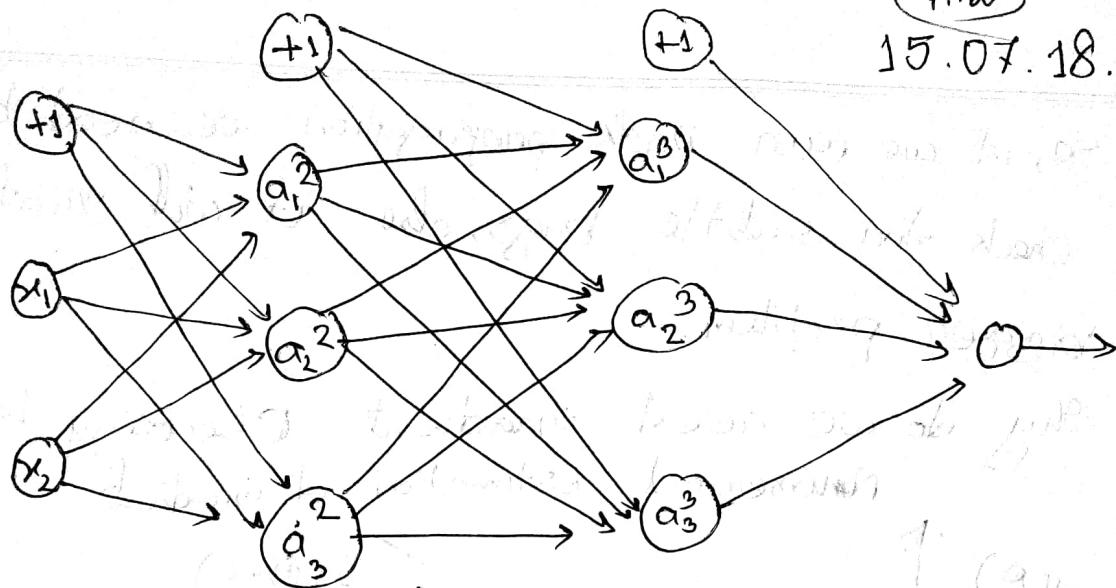
gradient base of

writing own stepsize

for each iteration

1) $\nabla J = (\theta) C \frac{\delta}{\epsilon \theta \delta}$

(Final)
15.07.18.



Gradient checking:

We know \rightarrow How to do forward propagation
 \rightarrow " " " backward

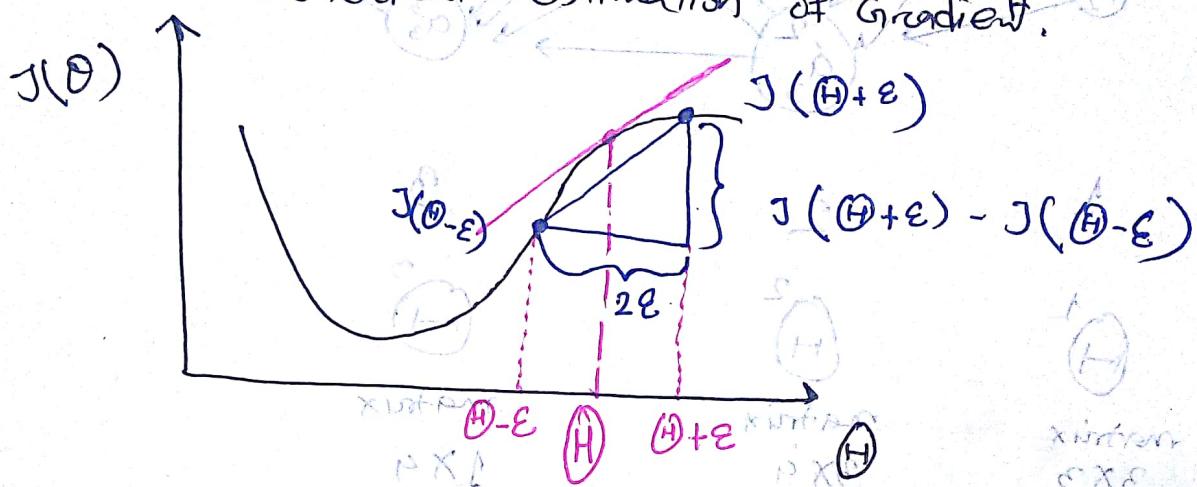
~~Backward~~ propagation

- has a lot of details
- tricky to implement
- there are many ways to have subtle bugs in back propagation.

So, if we run back-propagation, we need to check for subtle bugs, else it will create bigger problem.

* Why do we need Gradient Checking here?

Numerical estimation of Gradient.



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

$\tan \theta = \frac{dy}{dx}$

$$\frac{d}{d \theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \epsilon = 10^{-4}$$

Parameter Vectors

Unrolling Θ :

$$\Theta^1 = 10 \times 11$$

$$\Theta^2 = 10 \times 11$$

$$\Theta^3 = 11$$

$$\text{theta vectors} = [\text{theta1}(:), \text{theta2}(:), \text{theta3}(:)]$$

$$1-110 \quad 111-220 \quad 221-232$$

thetaVector[1]

Back to matrix,

$$\text{theta1} = \text{reshape}(\text{thetaVector}(1:110), 10, 11)$$

$$\text{theta2} = \text{reshape}(\text{thetaVector}(111:220), 10, 11)$$

Parameter vector Θ

Let $\Theta \in \mathbb{R}^n$ (Θ is unrolled version of

$$(\Theta^1, \Theta^2, \dots, \Theta^n)$$

$$\Theta = [\Theta_1, \Theta_2, \Theta_3, \dots, \Theta_n]$$

$$\frac{\partial}{\partial \Theta_{10}} J(\Theta)$$

$$D_{10} \quad \frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$D_{11} \quad \frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$$

(Grad approx. (1) & (2) both $\frac{\partial}{\partial \theta_i} J(\theta)$ is constant)

$$D_{1m} \quad \frac{\partial}{\partial \theta_m} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_m + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_m - \epsilon)}{2\epsilon}$$

(Grad approx. (3) $\frac{\partial}{\partial \theta_m} J(\theta)$ is constant)

for $i = 1 \dots m$

theta P = theta

theta QP(i) = thetaP(i) + EPSILON

theta M = theta

theta M(i) = thetaM(i) - EPSILON

grad Approx = $(J(\text{theta P}) - J(\text{theta M})) / 2 * EPSILON$

end

check that,

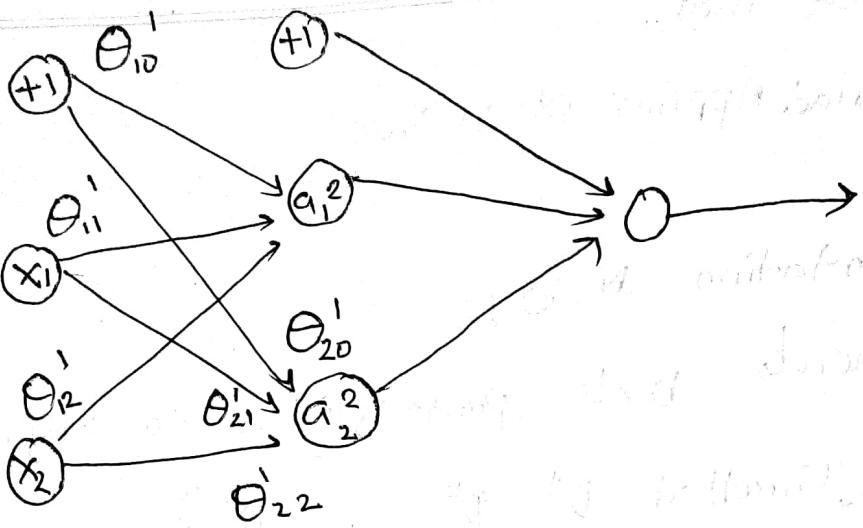
$$\text{grad Approx} \approx D \text{ Vec}$$

* Implementation Note:

1. Implement back propagation to compute D .
(Unrolled D^1, D^2, \dots, D^n)
2. Implement gradient check to compute grad approx.
3. Make sure they give similar result.
4. Turn off gradient checking, using back propagation code for learning, train your algorithm

Random Initialization:

- Initialize all the θ 's to 0 does not work with neural network.
- When back propagation all nodes will update to the same value repeatedly.
- Instead we can randomly initialize our θ 's for our Θ matrix.



$$a_1^2 = \theta_{10}^1 \theta_{11}^1 \theta_{12}^1 \text{ initialize}$$

$$\text{if } \theta_{10}^1 = 0 \text{ then } a_1^2 = 0 \text{ else calculate } a_1^2$$

$$a_2^2 = \theta_{20}^1 \theta_{21}^1 \theta_{22}^1 \text{ calculate the next } a_2^2$$

~~if $a_1^2 = 0$ then $a_2^2 = 0$ else calculate a_2^2~~

$$a_1^2, a_2^2$$

$$\Theta^T x \quad \Theta^T x$$

$$a_1^2 = a_2^2 \text{ and the multiplying signal } A$$

$$a_1^2 = a_2^2 \text{ and the multiplying signal } B$$

$$a_1^2 = a_2^2 \text{ and the multiplying signal } C$$

$$a_1^2 = a_2^2 \text{ and the multiplying signal } D$$

X-intercept \oplus two and

17.07.18.

Summary (all together):

First, pick a Network Architecture.

i.e.
3 3 4
3 4 4 4
3 5 5 5 4

↳ choose the layout of your Neural Network.

↳ # hidden units in each layer

↳ # layers in total Network.

- # input units = dimension of features ($x^{(i)}$)
- # output units = # classes
- # hidden units / layer = 

Second, Training a Neural Network:

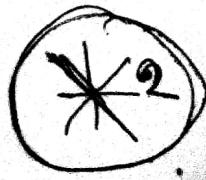
1. Randomly initialize the θ s.

2. Implement forward propagation to get $h_{\theta}(x)$ for any $x^{(i)}$

3. Implement the cost function $J(\theta)$

4. Implement Back propagation to compute

partial derivative. $\frac{\partial}{\partial \Theta_{ij}} J(\theta)$



Symmetry Breaking & Random Initialization.

Initial each $\Theta_{ij}^{(l)}$ to random value in

$$[-\epsilon, \epsilon]$$

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



Hidden unit / layer = Usually more is better.

must balance with cost as it increases with
more Hidden Unit.

• Default : 1 hidden layer

If you have more than one hidden layer,
you should have same Hidden
units in every hidden layers. (Suggested)



When perform Back propagation & forward propagation, we loop every training example!

for $i = 1 : m$.
Perform FP & BP using $(x^{(i)}, y^{(i)})$
(Get activations $a^{(l)}$, delta terms $\delta^{(l)}$ for $l = 2, 3, \dots, L$)

$\downarrow \text{for}$

Then compute $D^{(t)}$ out of for loop.

5 Use gradient checking to compare $\frac{\partial}{\partial \theta_{ij}} J(\theta)$.

Compute using BP vs using numerical estimation of gradient of $J(\theta)$ to confirm that your BP works. Then disable gradient checking.

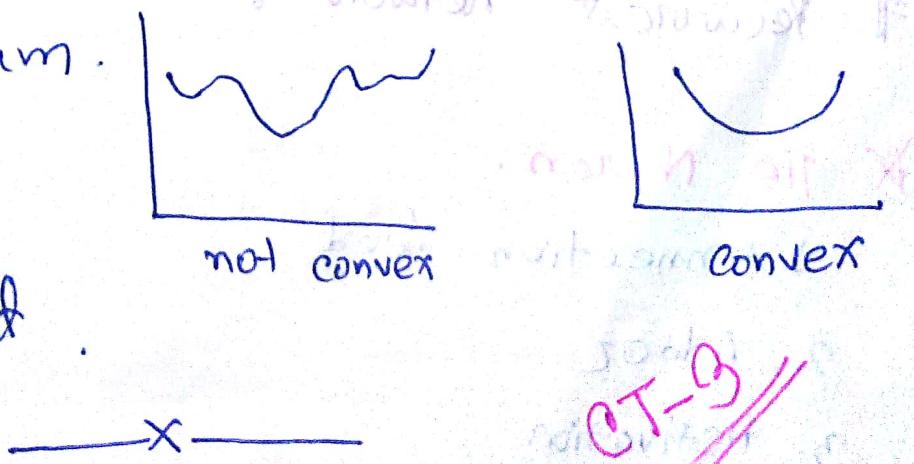
6. They use GD or a "builtin" optimization function to minimize the cost function with weight in theta.

→ $J(\theta)$ in NN is not a convex function.

It could be local minimum as well as global minimum.

↓
a good local

optimal can get.



"Special class"

19.07.18.

Neural Networks:

Zain Protein Molecule,

Human's Unsupervised learning \rightarrow Imagination

Architecture of neuron or Cell to discuss

1. Single-layer feed-forward

2. multi-layer "feed" or "deep" learning

3. Recurrent

Q: How to choose hidden layers' number
node?

Data prepossessing \rightarrow fit data in model

Recurrent network:

* The Neuron.

1. Connection feed

2. Adder

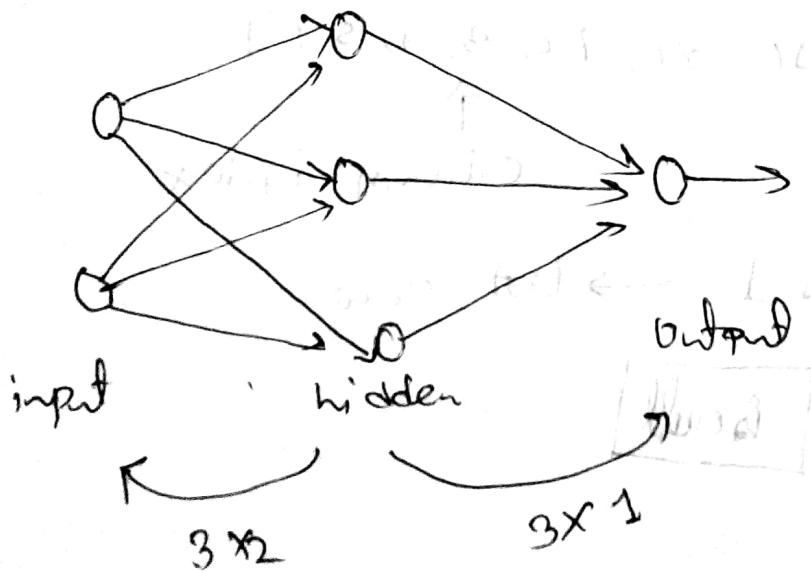
3. Activation

Dimensions of a NN:

1. Neurons
2. Architecture
3. Learning algo.
4. Applications

* learning Algorithm: Backpropagation:

$$w'_{46} = w_{46} + \eta \delta \frac{df_6(e)}{de} y_4$$



2.2.2.1

"Python"

Numpy → Numerical python

Scipy → Scientific python

Pandas

SciKit - learn

matplotlib

Seaborn

Seaborn.

Pandas → all data processing

[10 : 20, [0, 3, 4, 5]]



column index

df.iloc[-1] → last row

isnull

Linear regression.

epoch size

Association Rule Mining:

Market Basket Analysis

1. Identify frequent item set \rightarrow based on "support."
2. Based on frequent item set, discover association rule \rightarrow satisfying "confidence."

i.e.: Consider the following transactions:

<u>t₁</u>	Bread, Butter, Milk
<u>t₂</u>	Bread, Butter, Coke, Diaper
<u>t₃</u>	Coke, Diaper, Butter
<u>t₄</u>	Bread, Milk, Coke, Diaper
<u>t₅</u>	Coke, Butter, Milk

Requirement: Minimum Support = 70% (of total transaction)
(Support) Minimum Confidence = 70%

Single item

$$\{\text{Bread}\} = \frac{3}{5} = 0.6 \times$$

$$\{\text{Butter}\} = \frac{4}{5} = 0.8 \checkmark$$

$$\{ \text{milk} \} = \frac{3}{5} = 0.6 \quad \times$$

$$\{ \text{Diaper} \} = \frac{3}{5} = 0.6 \quad \times$$

$$\{ \text{Coke} \} = \frac{4}{5} = 0.8 \quad \checkmark$$

* {Butter}, {coke}

Two items

$$\{ \text{Butter, coke} \} = \frac{3}{5} = 0.6 \quad \times$$

Association Rule (confidence)

$$\{ \text{Butter} \rightarrow \text{Coke} \} = \frac{3}{4} = 0.75 \quad (75\% \text{ chance})$$

↑
when buys butter,
then buys coke

(Apriori Algorithm) → (FP Growth Algorithm)

(Best → frequent pattern (FP))

How many Butter, coke (Faster)

How many Butter

Min support & min confidence
depends on customer requirements.

$$0.8 = \frac{3}{4} = \{ \text{Butter} \}$$

could be

{ Butter \rightarrow Diaper, coke }

* All possible combination

{ coke \rightarrow Butter, Diaper }

depends on

{ Diaper \rightarrow Butter, coke }

min. support

Butter, Diaper, coke

Butter

"Decision tree"

Test Data

* Decision tree classification:

↓
Supervised learning

Phases

Input training set

(Training Data)

↓
Inference

to model creation

Decision Tree algo.

model

(Decision tree
model)

classify test data
(output)

Blast

i.e. Example training Data:

SL	Color	Age	Class
1	white	15	Tender
2	white	20	Tender
3	Gray	25	Mid range
4	Gray	30	Mid range
5	Red	35	Mature
6	Red	36	Mature

1. calculate "Entropy"

2. calculate "Information gain" for each attribute except the class.

* Entropy:

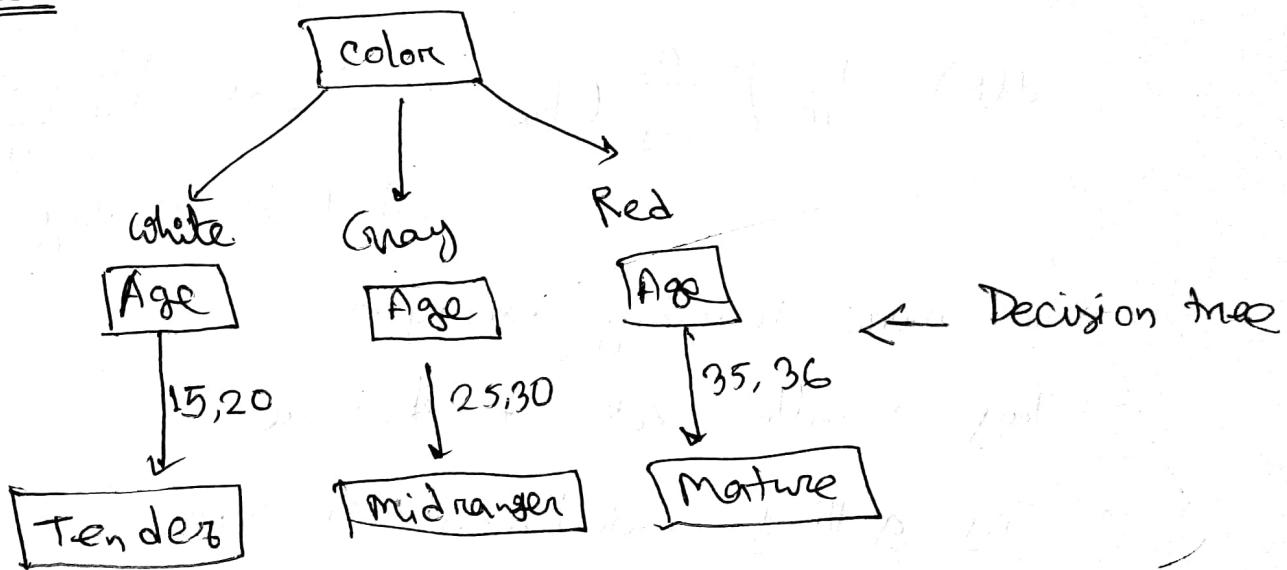
→ Inconsistency in data. [How entropy is expected]

→ Higher the entropy, the higher the inconsistency.

Information Gain

→ To decide which attribute to use as "root" of the tree.

Trees



New Data:

Red 32 (Not available).

— X —

22.07.18.

- ML system Design
 - Understand multiple parts.
 - How to deal with skew data.
- * Regularized Linear Regression (LR)

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

what should you try next?

- Get more training examples.
- Try smaller set of features.
- Try getting additional features.
- Try adding polynomial features
 $(x_1^2, x_2^2, x_1 x_2, \text{etc})$
- Decrease λ
- Increase λ

Do not try anything randomly

ML Diagnostic:

A test that you can run to gain insight what is/ is not working with a learning algorithm, and gain guidance as to how best to improve its performance.

Evaluating Hypothesis: (HT):

→ A HT may have lower error for the training examples, but still be inaccurate.
(i.e. overfitting problem)

→ Thus evaluate a HT, given a data set up of training example, we can split the data into two sets,

- ① Training set: (70% Data may be)
1st ↗
- ② Test set: (30% ")
last ↗

x^1, y^1
 x^2, y^2
 \vdots
 x^m, y^m

Training Set (70%)

$x^1_{\text{test}}, y^1_{\text{test}}$
 $x^2_{\text{test}}, y^2_{\text{test}}$
 \vdots
 $x^m_{\text{test}}, y^m_{\text{test}}$

Test set (30%)

- New procedure using these two set is then
1. learn training Θ and minimize $J_{\text{train}}(\Theta)$ using test set.
 2. Compute the test set errors $J_{\text{test}}(\Theta)$
↳ LR

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

\hookrightarrow classification Problem:

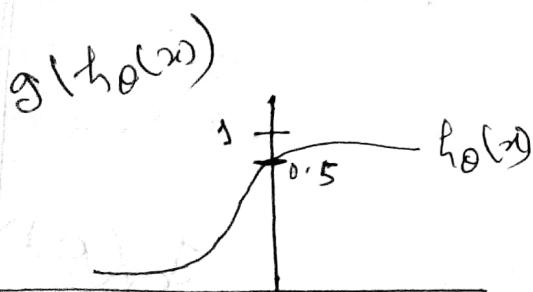
$$J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left[y_{\text{test}}^{(i)} \log(h_{\theta}(x_{\text{test}}^{(i)})) + (1 - y_{\text{test}}^{(i)}) \log(1 - h_{\theta}(x_{\text{test}}^{(i)})) \right]$$

Missclassification error,

$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5 \text{ &} \\ & y = 0 \\ 0 & \text{or, } h_{\theta}(x) < 0.5 \text{ &} \\ & y = 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}}$$

$$\sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$



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

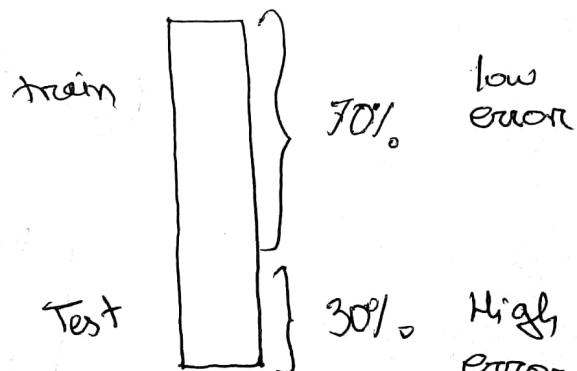
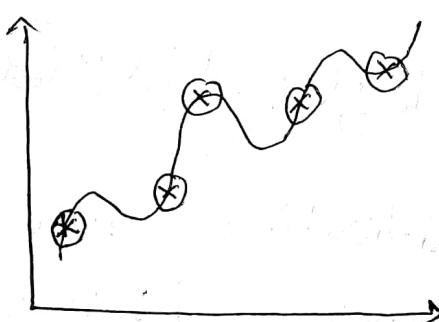
24.07.18.

$$\theta^T x \quad \theta x$$

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_3$$

$$J(\theta) \rightarrow \min$$



(as not trained
with these data)

Model Selection:

$$d=1 \quad 1. \quad h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow J_{\text{train}}(\theta^1) \rightarrow \theta^1 \rightarrow J_{\text{test}}(\theta^1)$$

$$d=2 \quad 2. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow J_{\text{train}}(\theta^2) \rightarrow \theta^2 \rightarrow J_{\text{test}}(\theta^2)$$

$$d=3 \quad 3. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \rightarrow J_{\text{train}}(\theta^3) \rightarrow \theta^3 \rightarrow J_{\text{test}}(\theta^3)$$

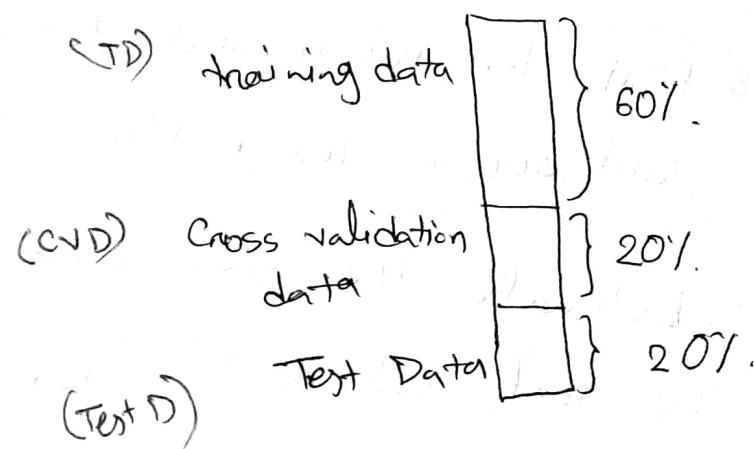
...

$$d=10 \quad 10. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{10} x^{10} \rightarrow J_{\text{train}}(\theta^{10}) \rightarrow \theta^{10} \rightarrow J_{\text{test}}(\theta^{10})$$

$d = \text{degree of polynomial}$

When, $d = 5$, $J_{\text{test}}(\theta^5)$ [assume] will give lowest value.

* d has been trained using test data (30%).



* Break down data set into three sets:

- training set : 60%
- Cross validation set : 20%
- Test set : 20%

* We can now calculate three separate error values for the three different sets using the following method.

1. Optimize the parameters in \hat{H} using the training set for each polynomial degree.

2. Find the polynomial degree d with the least error using the cross validation set.

3. Estimate the generalization error using the test set with $J_{\text{test}}(\theta)$.

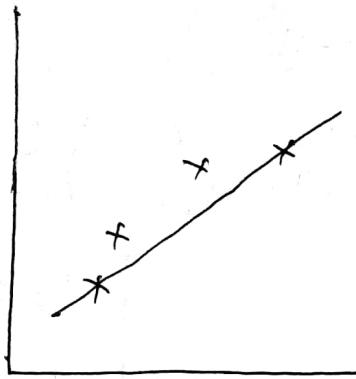
($d = \text{theta from polynomial with lower error}$)

* This way the degree of the polynomial d has not been trained using the test set.



29.07.18.

Diagnosing Bias Vs Variance:

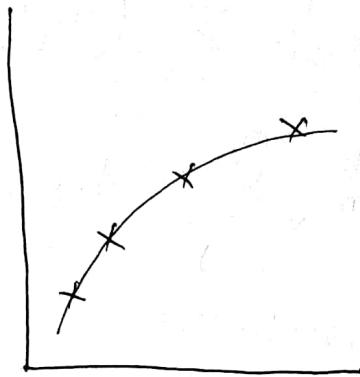


$$\theta_0 + \theta_1 x$$

High Bias

Underfit

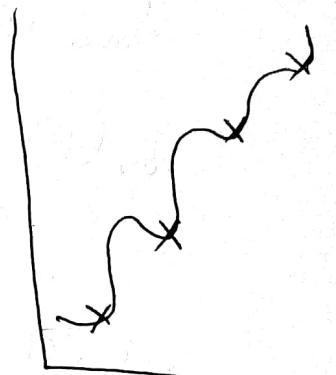
$$d = 1$$



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

Just right

$$d = 2$$



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

overfit

High Variance

$$d = 4$$

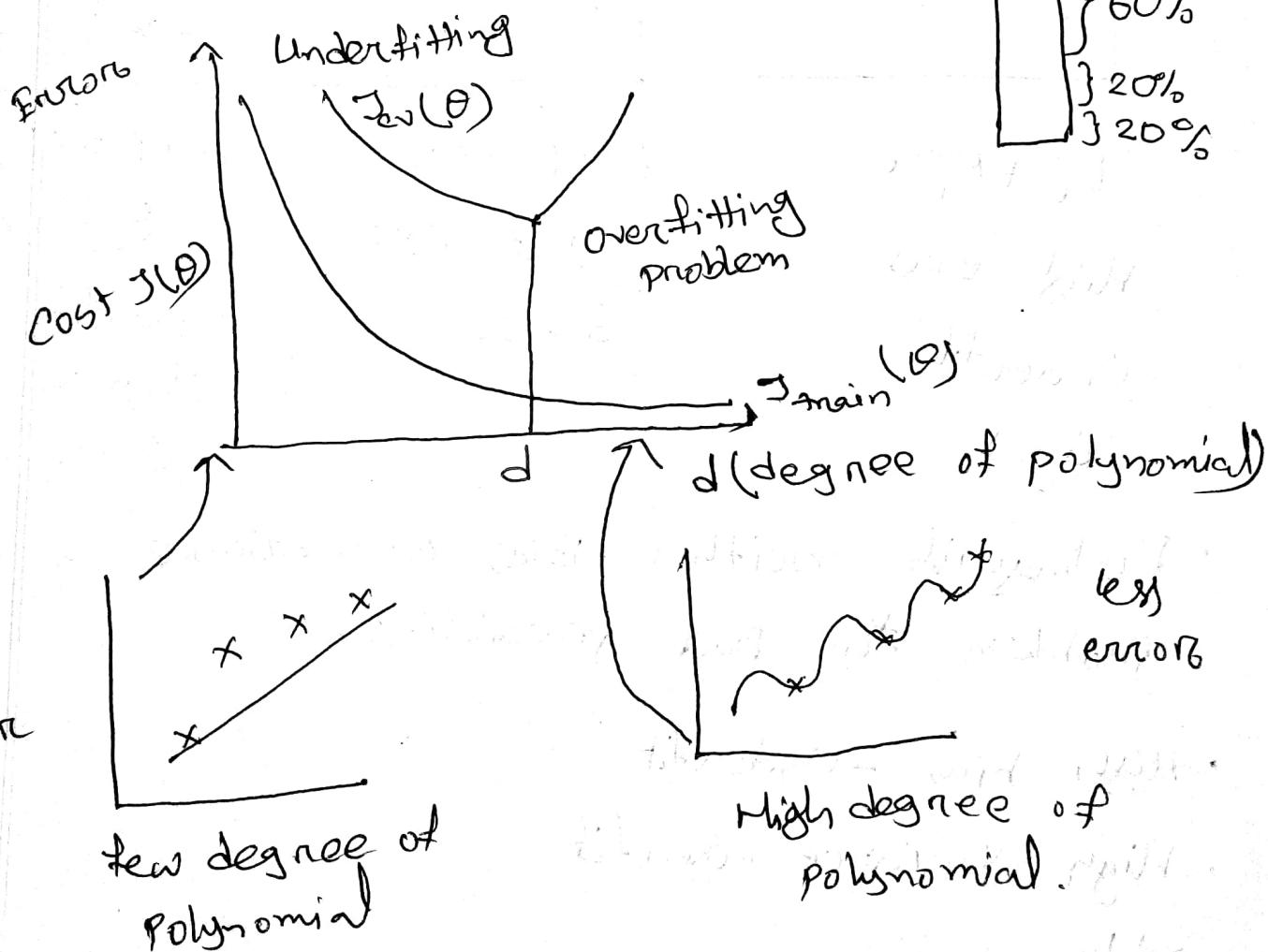
- Distinguish whether Bias or variance is the problem for Bad prediction.
- High bias - Underfit
- High variance - Overfit
- golden ratio.

cv = cross validation.

Bias/variance:

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

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

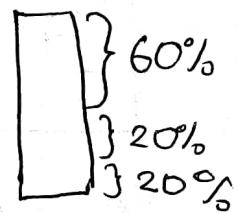
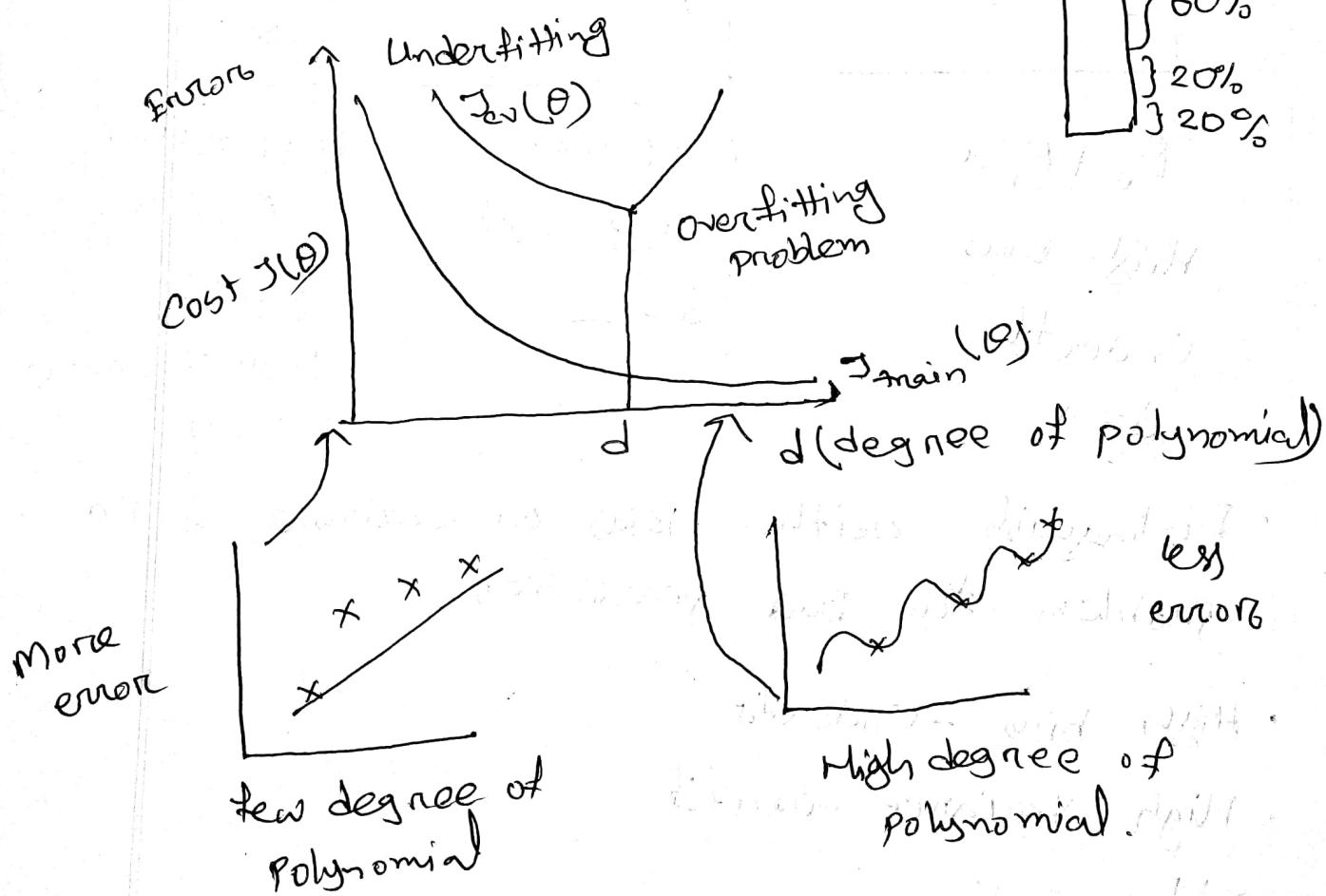


cv = Cross validation.

Bias/variance:

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

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



For high Bias,

$J_{\text{train}}(\theta)$ & $J_{\text{cv}}(\theta)$ will be high

$$J_{\text{train}}(\theta) \approx J_{\text{cv}}(\theta).$$

For high variance,

$J_{\text{train}}(\theta) \rightarrow \text{low}$ and $J_{\text{cv}}(\theta) \rightarrow \text{high}$

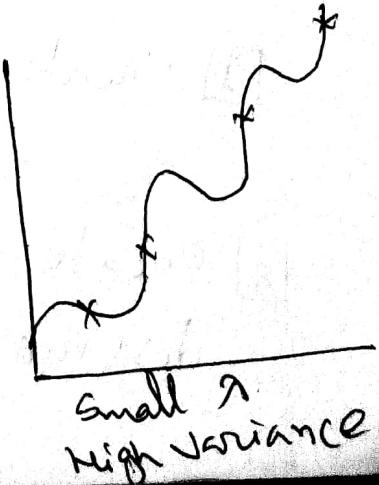
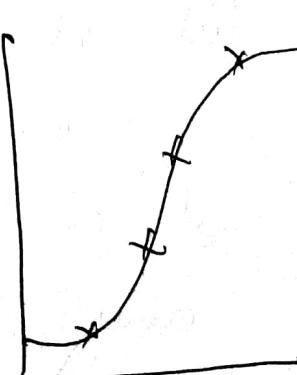
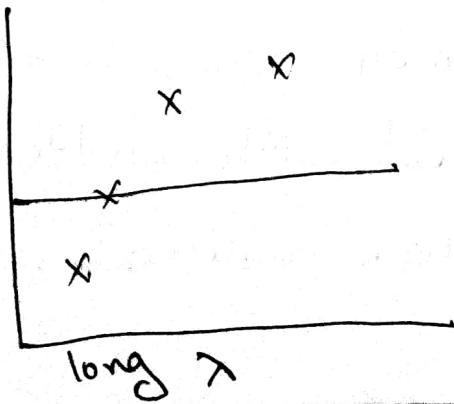
$$J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta).$$

Regularization & Bias/Variance:

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

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

$$J_{\text{cv}}(\theta) = \frac{1}{2M_{\text{cv}}} \sum_{i=1}^{M_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - \hat{y}_{\text{cv}})^2.$$



$$\theta_0 \approx 0, \theta_1 = 0 \dots \theta_n \approx 0$$

$$h_{\theta}(x) \approx \theta_0$$

After selecting d, How to Judge ?

1. Try $\lambda = 0 \rightarrow \min J(\theta) \rightarrow \textcircled{H}^{(1)} \rightarrow J_{cv}(\theta)$

without regularization

2. $\lambda = 0.01 \rightarrow \min J(\theta') \rightarrow \textcircled{H}^{(2)} \rightarrow J_{cv}(\theta^2)$

$\lambda = 0.02 \rightarrow \textcircled{H}^{(3)} \rightarrow J_{cv}(\theta^3)$

$\lambda = 0.04 \rightarrow \textcircled{H}^{(4)} \rightarrow J_{cv}(\theta^4)$

$\lambda = 10.24 \rightarrow \min J(\theta) \rightarrow \textcircled{H}^{(10)} \rightarrow J_{cv}(\theta^{10})$

λ , with less error will be selected.

Choose the model and λ :

1. Create a list of λ 's (i.e. $\lambda \in \{0, 0.01, 0.02, 0.04, \dots, 10.24\}$)

2. Create a set of model with different degrees or any other variants.

③ Iterate through the λ 's and for each
 λ go through all models to learn Θ .

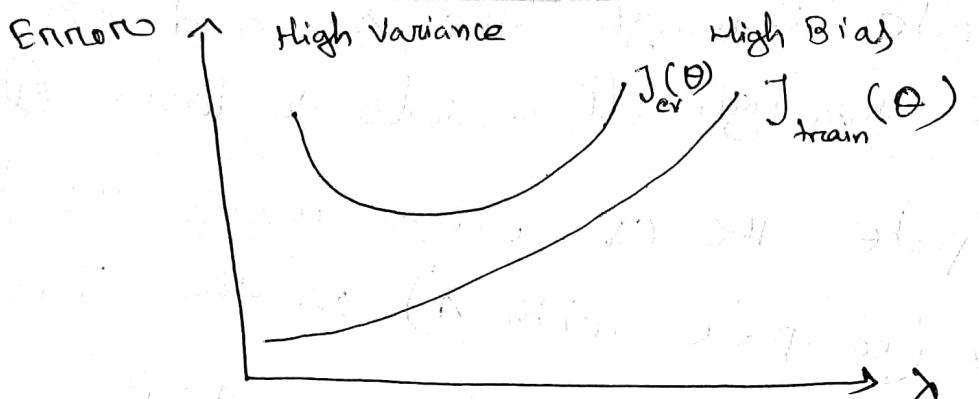
④ Compute the CV error using the learned
 Θ (compute with λ) on the $J_{cv}(\theta)$
without regularization or $\lambda = 0$

⑤ Select the best combination that produces
the lowest error on the cross validation
set.

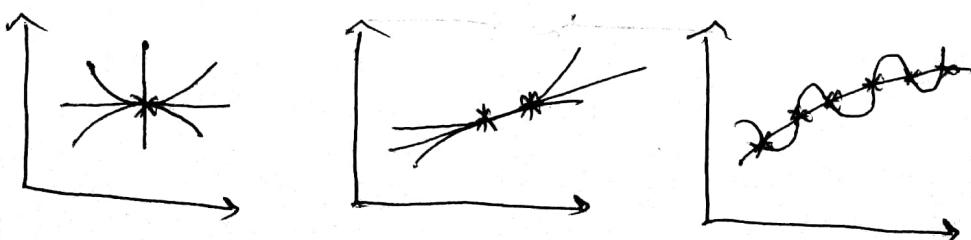
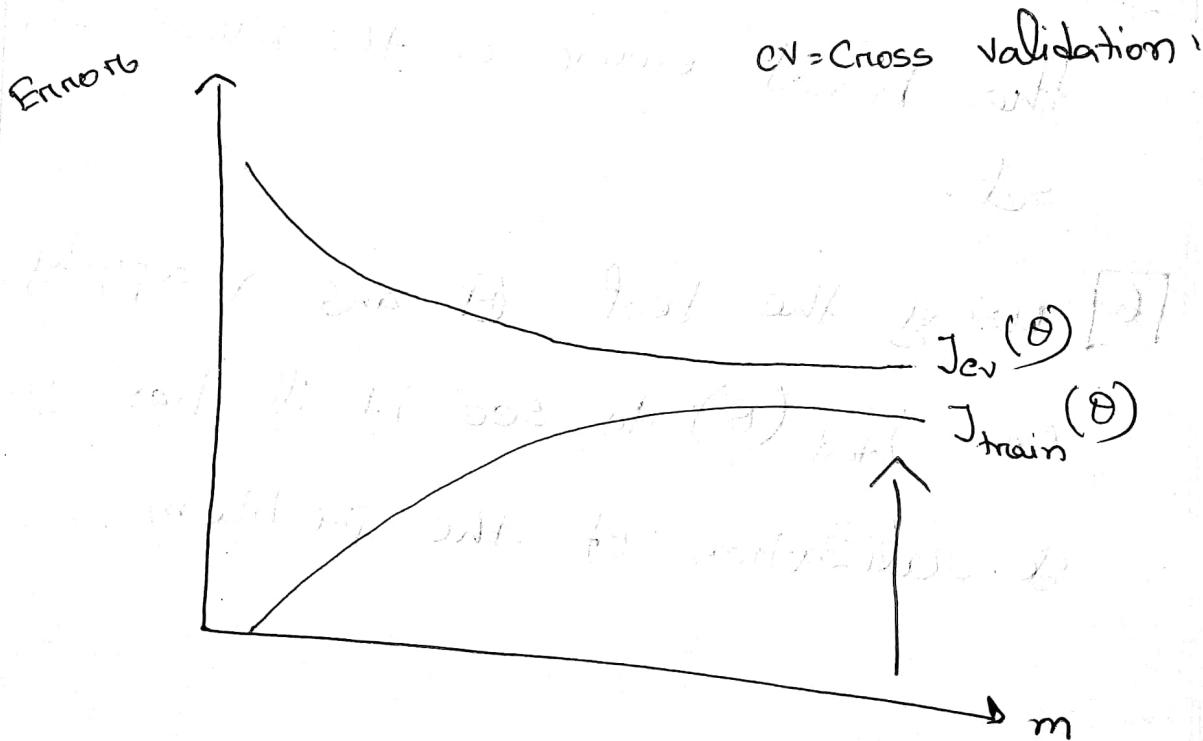
⑥ Using the best Θ and λ apply it
on $J_{test}(\theta)$ to see if it has good
generalization of the problem.

— \rightarrow —

31.07.18



Learning Curves: Bias and Variance



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

→ Train algo. on a very few data, which

- ↳ easily have 0 error.
- ↳ can find a quadratic function that touches all point.

→ Training set larger,

- ↳ Error increases

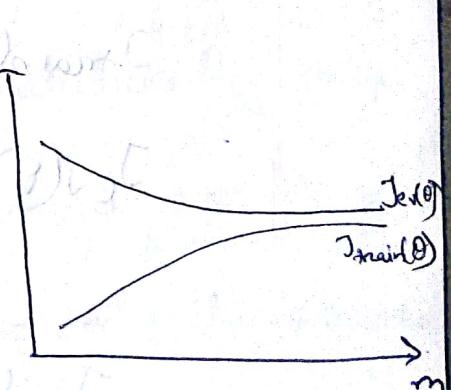
→ Error value ^{will} plateau out after a certain m.

Experiencing High Bias:

→ Lower training set size,

$J_{train}(\theta) \rightarrow$ low

$J_{cv}(\theta) \rightarrow$ high



→ Higher training set size,

$J_{train}(\theta) \rightarrow$ high

$J_{cv}(\theta) \rightarrow$ high

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

* If learning algos. is suffering from high bias, getting more data will not (by itself) help much.

Experiencing High Variance:

→ Lower training set size,

$$J_{\text{train}}(\theta) \rightarrow \text{low}$$

$$J_{\text{cv}}(\theta) \rightarrow \text{High}$$

→ Higher training set size,

$$J_{\text{train}}(\theta) \rightarrow \text{increase}$$

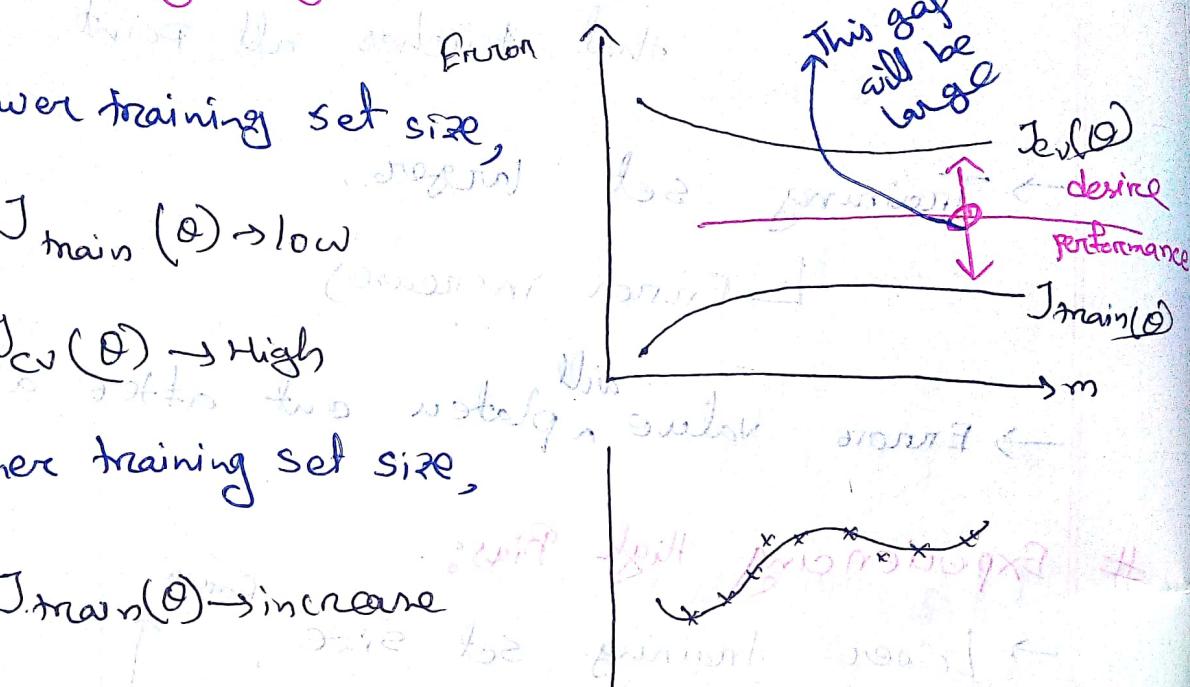
$$J_{\text{cv}}(\theta) \rightarrow \text{Decrease}$$

(But not low)

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

$J_{\text{cv}}(\theta) \rightarrow$ Decrease without leveling off.

* difference must be significant.



Colab

λ = regularization parameter

- * What to do next? How to fix?:

- Getting more training examples : Fixes HV
- Trying smaller features : Fixes HV
- Adding features : Fixes HB
- Adding polynomial features : Fixes HB
- Decrease λ : Fixes HB
- Increase λ : Fixes HV

Diagnosing Neural Network:

- Fewer Parameters.
 - Prone to underfitting
 - Computationally cheaper.
- Higher Parameters.
 - Prone to overfitting (λ increase)
 - Computationally expensive.

- * Using a single hidden layer is a good default.
You can train yours of hidden layers using CV set.
You can then select the one that perform most.

28.08.18.

Handling Skewed Data:

↳ Have a lot more of example from one class than from the other class(es).

* Cancer classification Problem:

Train logistic model $h_0(x) = \begin{cases} 1 & \text{if cancer} \\ 0 & \text{if not} \end{cases}$

Your algorithm got 99% accuracy

0% Error.

0.05% Error

function $\text{PnC}(x)$

$y=0$ # ignore x

99.5% Accuracy

0.05% Error

Error Matrix:

$Y=1$. cancer

Precision / Recall

$\log(x) = 1$

Actual class

Predicted class

Actual class	
1	0
True Positive	True Positive
	False Positive
False Negative	True Negative

* Precision: of all patients where we predict $Y=1$, what fraction actually has cancer?

$$P = \frac{\text{True Positive}}{\text{Total no. of Predicted Positive}}$$

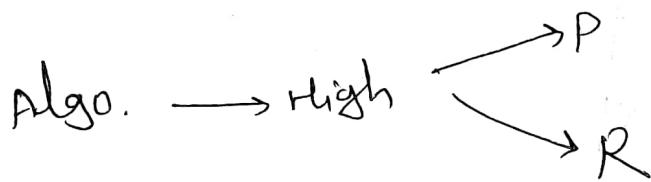
$$= \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

* Recall: of all patient who have cancer what fraction did we correctly detect as having cancer?

* High precision & Recall gives better performance

$$R = \frac{\text{True Positive}}{\# \text{ Actually Positive}}$$

$$= \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$



* It is not possible to get high P and R at the same time setting $Y=1$ or $Y=0$ always.

If $Y=1$, $P \rightarrow 0$

If $Y=0$, $R \rightarrow 0$

P and R

Logistic Regression

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

Predict 1, if $h_{\theta}(x) \geq 0$

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

Predict $Y=1$, only when very confident

Predict, $Y=1$ if $h_{\theta}(x) \geq 0.7$

" $Y=0$ if $h_{\theta}(x) < 0.7$

(Higher P, Lower R)

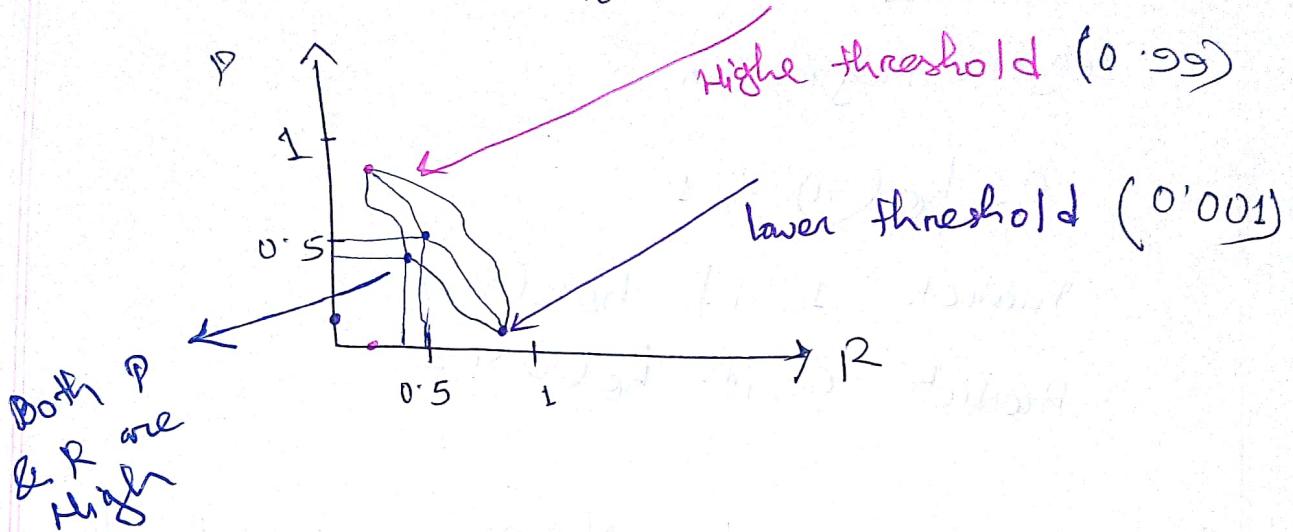
* Want to avoid missing.

Predict 1, if $h_{\theta}(x) \geq 0.3$

" 0, if $h_{\theta}(x) < 0.3$

(Lower P, Higher R)

Predict 1 if $h_0(x) > \text{threshold}$.



F₁ Score (F score)%

	$\frac{P}{R}$	$\frac{(P+R)/2}{F_1}$
Algo. 1	$0.5 / 0.4$	$0.45 / 0.44$
Algo. 2	$0.7 / 0.1$	$0.4 / 0.175$
Algo. 3	$0.02 / 1.0$	$0.50 / 0.0392$

If $Y = 1$ (always) $R = 1$ ~~P ≈ 0~~ $P \approx 0$

If $Y = 0$ $R = 0$ ~~P ≈ 1~~ $P \approx 1$

$$F_1 = \frac{2PR}{P+R}$$

* what will be the approach in case of choosing algo of large Data? (v.v.g)

$$F_1 = \frac{2PR}{P+R}$$

if $P = 0$ $F_1 = 0$

$R = 0$ $F_1 = 0$

if $P = 1$ $R = 1$ then $F_1 = 1$

* What is Support Vector Machine? [Search]
(SVM)

— X —

Support Vector Machine:

04.09.18

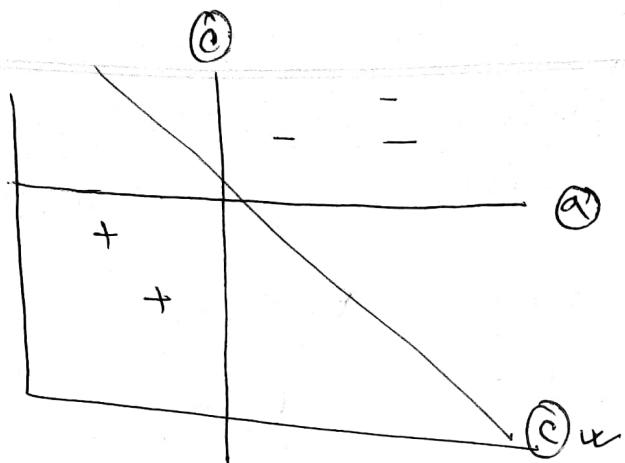
"It's not who has the best algo. that wins.

It's who has the most data."

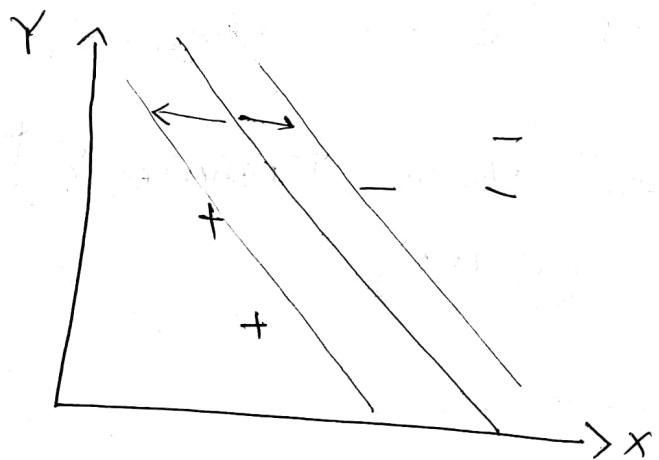
When TRUE, when FALSE?

i.e. For breakfast I ate one egg.

to, two, too ?



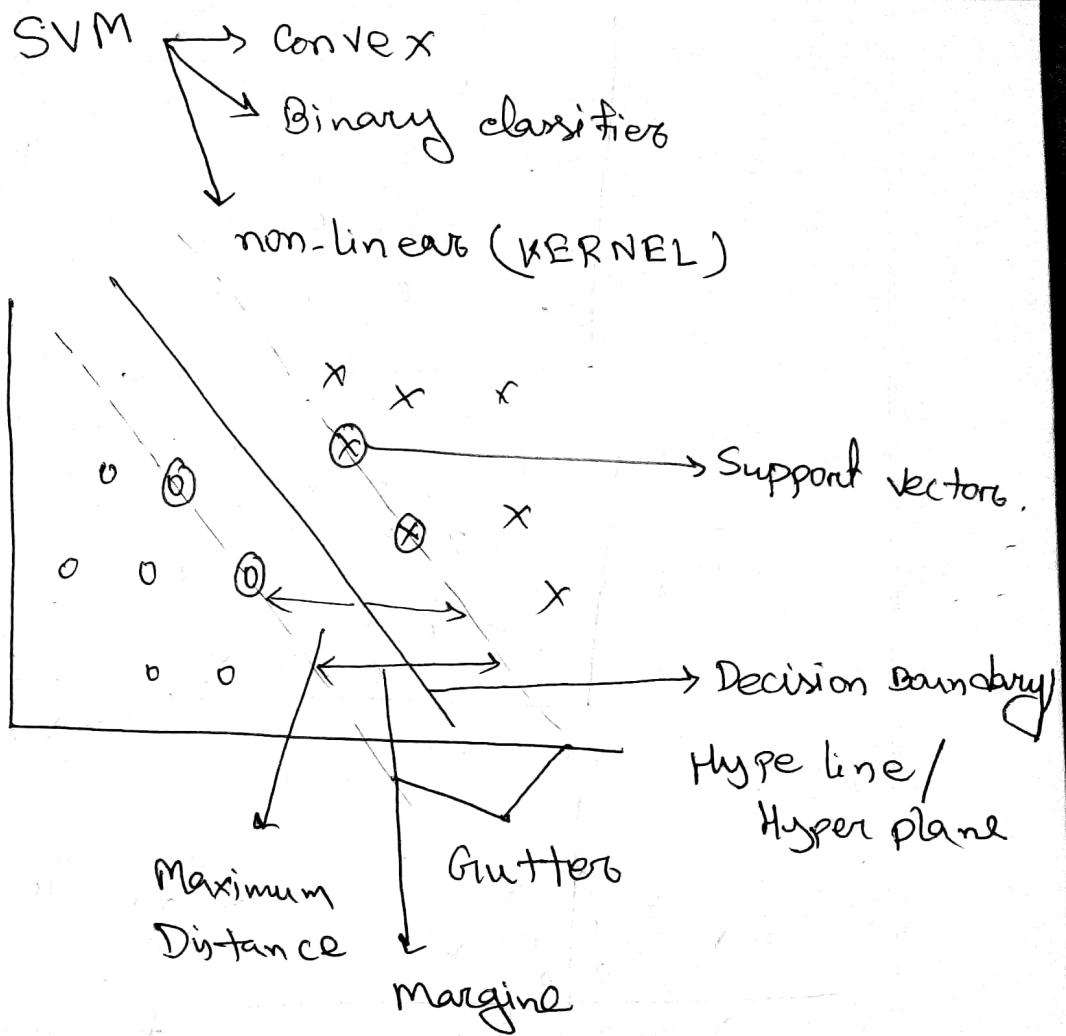
widest street
approach.



SVM \rightarrow ML. Algo. (Supervised)

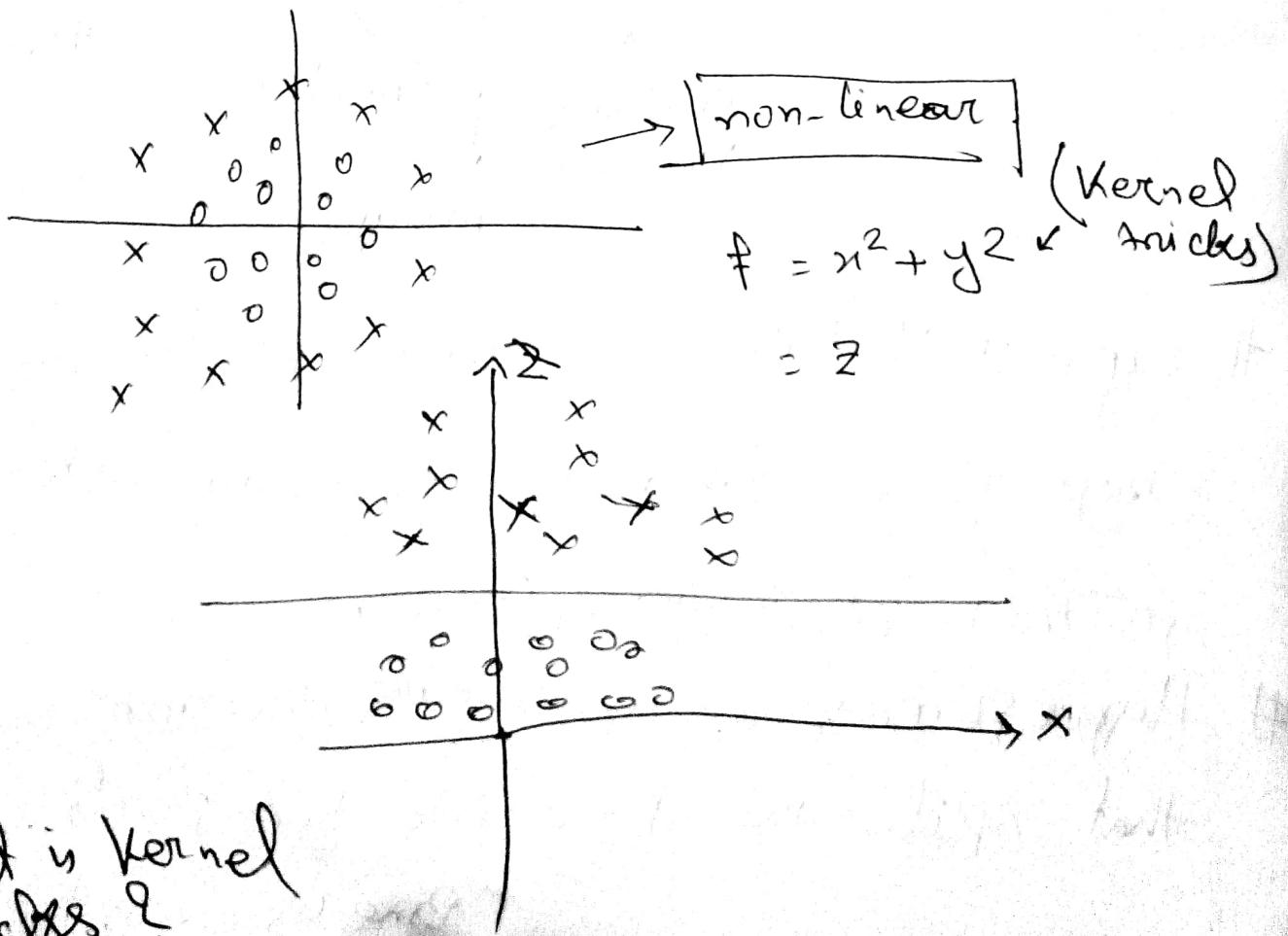
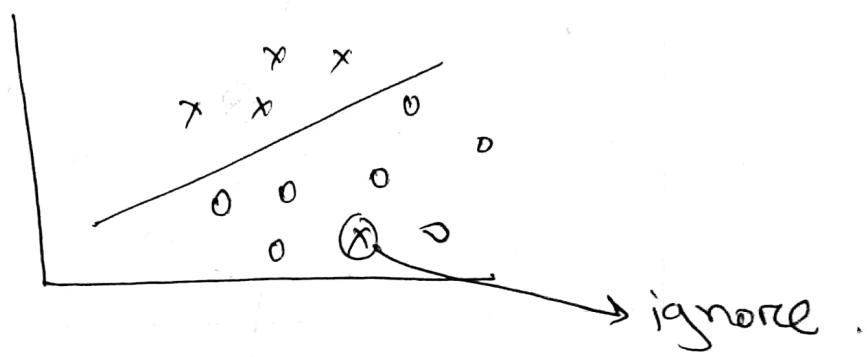
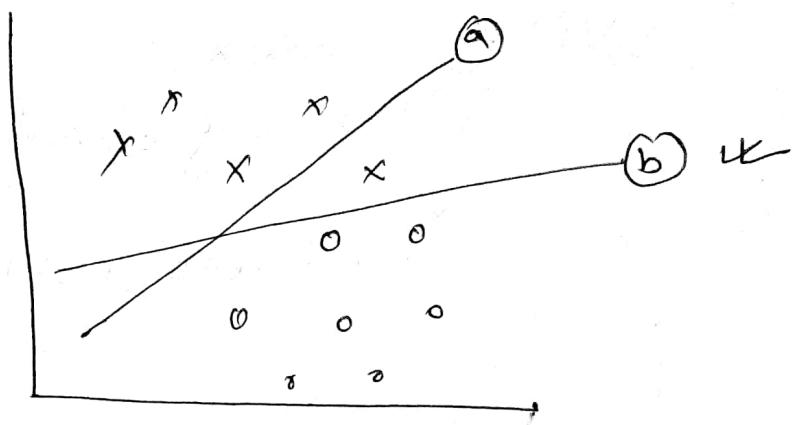
\searrow Numerical classifiers
regression/classifiers,

- \Rightarrow That draws a single decision boundary
to maximize the margin between
two region.



- # Support Vectors are data points nearest to hyper plane. If removed, would alter the position of hyper plane.
- # Hyper planes is a linear decision surface that splits the space into two parts.

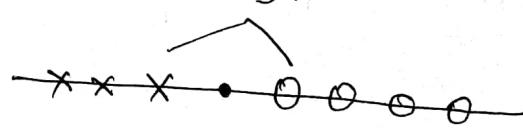
Binary classifiers



* What is Kernel tricks?

* # of Support Vectors (SV) :

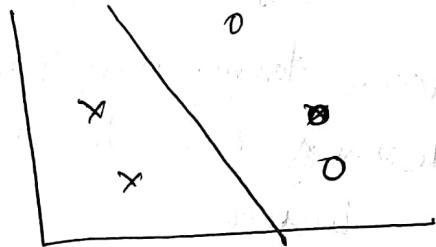
1 Dimension :



of SV = 2

D : point.

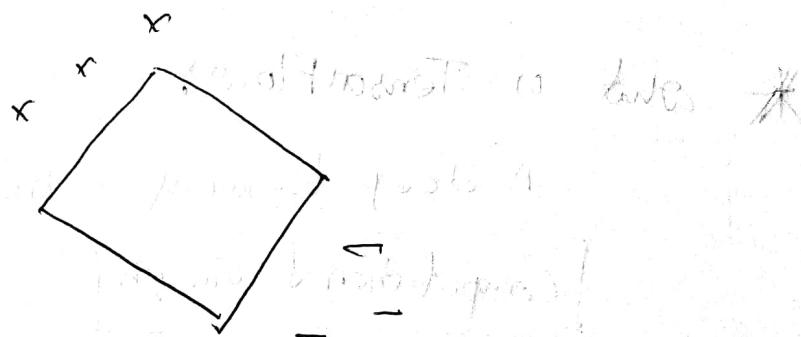
2 Dimension :



of SV = 2/3/4

line.

3 Dimension :



* How to select SV? [Convex rule]

→ H.G
Jessica Notes [MIT]
(AI)
2016

"Special Class"

06.09.18.

Wearable Computing

Ubicomp

Machine Learning using TensorFlow library:

Machine learning Vs Deep learning (DL)

learn from Solution → DL
Learning over learning
layers

DL → representation learning
Hierarchical representation
Dynamic learning

* What is TensorFlow?

A deep learning library

Computational Graph

import tensorflow as tf

Machine translation flow to TensorFlow.

```
import tensorflow as tf
```

```
a = tf.constant(1)
```

```
b = tf.constant(2)
```

```
c = tf.add(a, b)
```

```
session = tf.Session()
```

```
value_of_c = session.run(c)
```

Graph building

prepare execution env

initialize variables

Run the computation.

Two computation Phrases of a graph:

(i) construction phrase

(ii) Execution Phrase

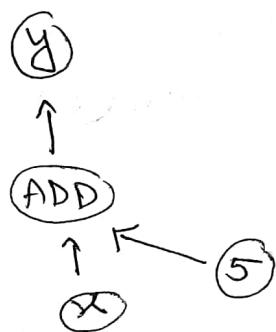
Everything is an object in TensorFlow.

* What is a Tensor?

Tensor data structure to represent all data.

* Basic of TensorFlow:

① Variable



* linear Algebra:

① Dimensionality

(1)

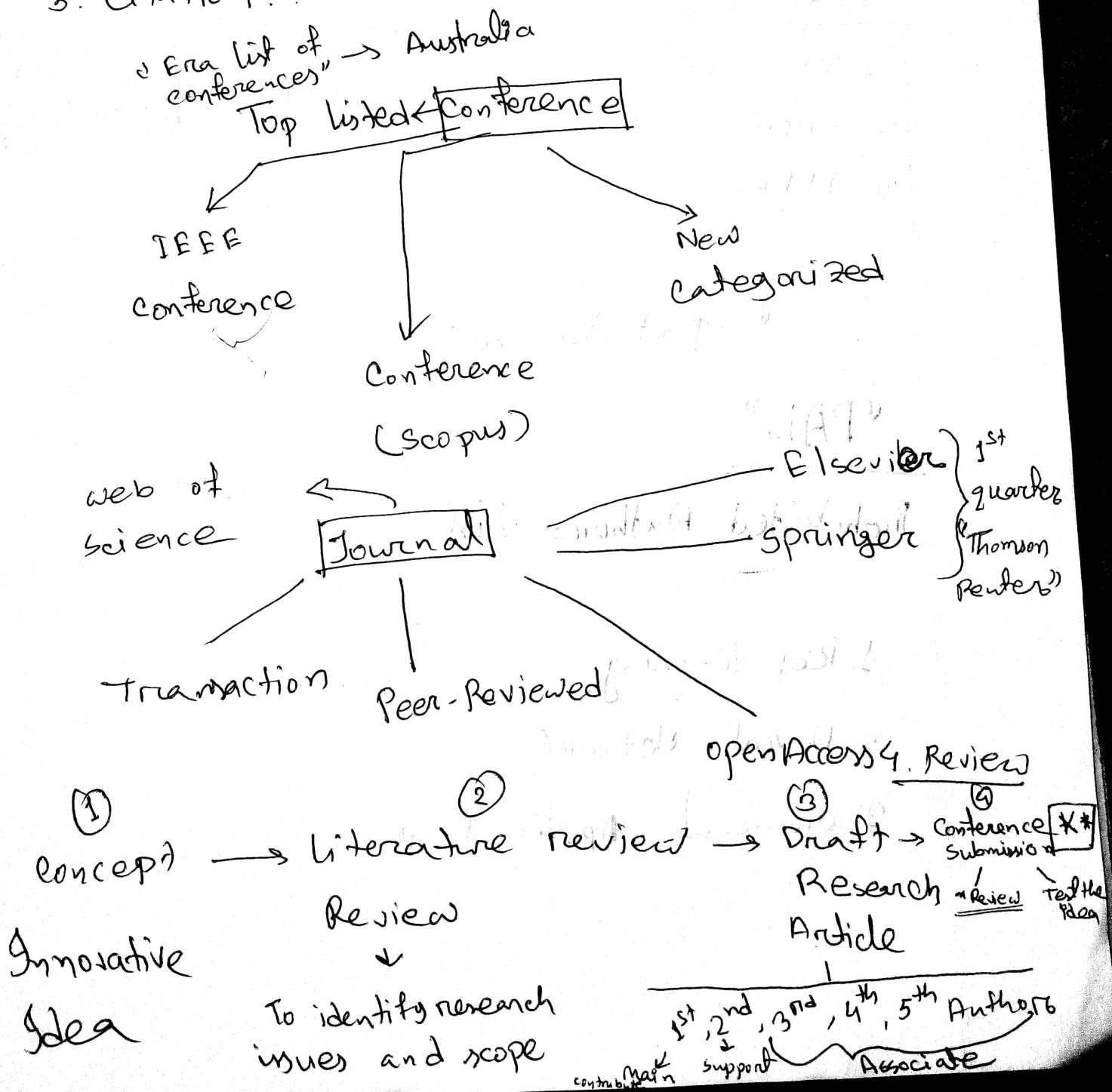
$$\begin{matrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{matrix}$$

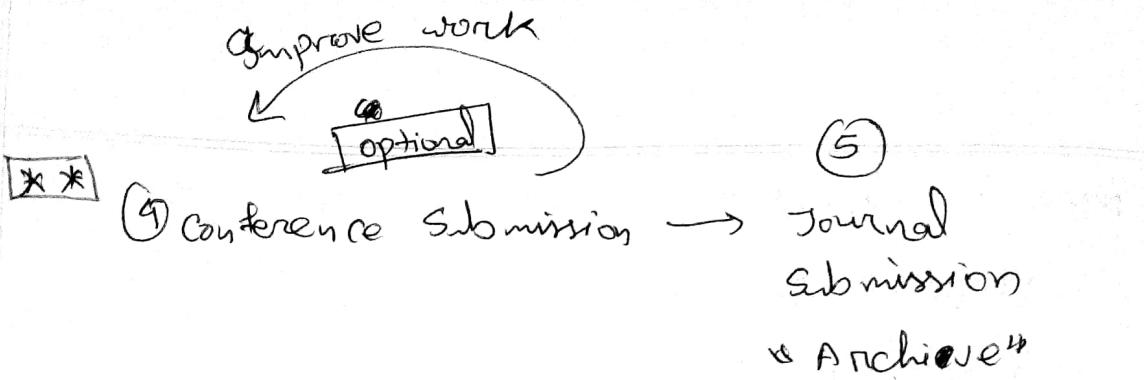
$$\begin{matrix} 2 & 4 & 6 \\ 8 & 10 & 12 \end{matrix}$$

linear Regression:

Research Discussion:

1. Plagiarism
2. Paraphrasing
3. Citation.





Indexing

1. ISI Index
2. Scopus
3. IEEE

"Thomson Reuters"

"impact factor"

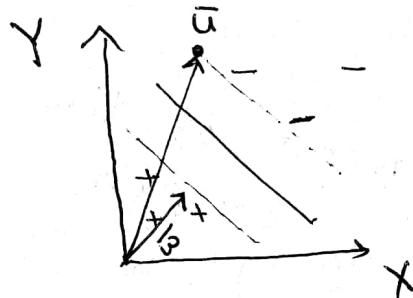
"PAL"

Prohibited Authors List

1. Deep learning
2. Neural Networks
3. General Applications

09.09.18.

Support Vector Machine:



Vector is perpendicular
with my hyperplane,

u = unknown point / new
example.

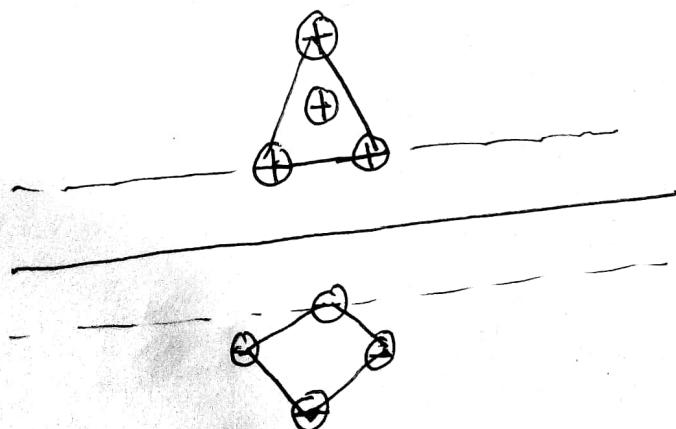
$$\text{if, } \bar{w} \cdot \bar{u} \geq c$$

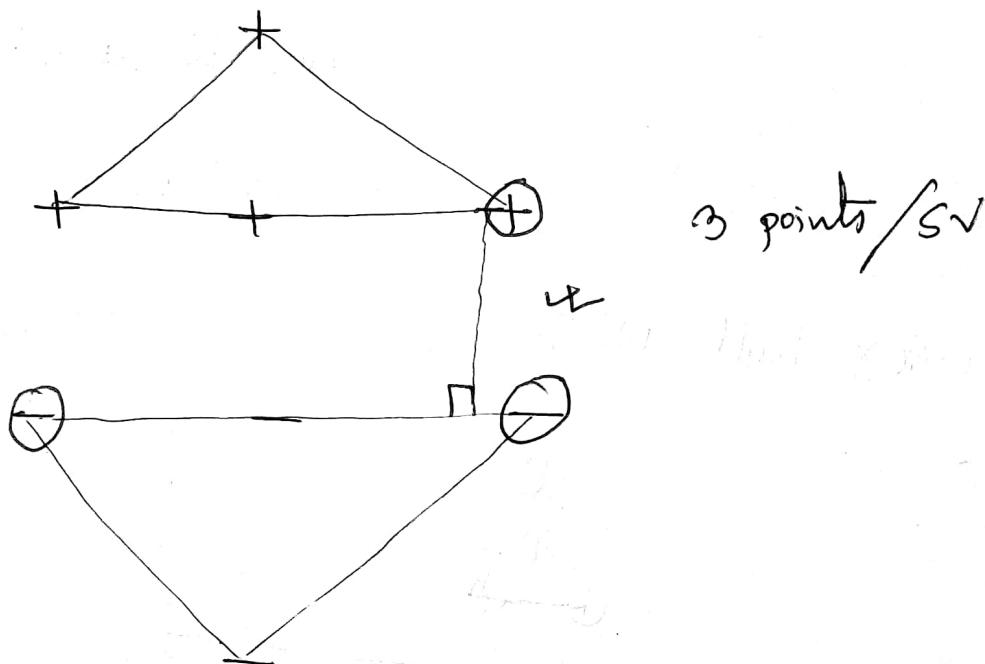
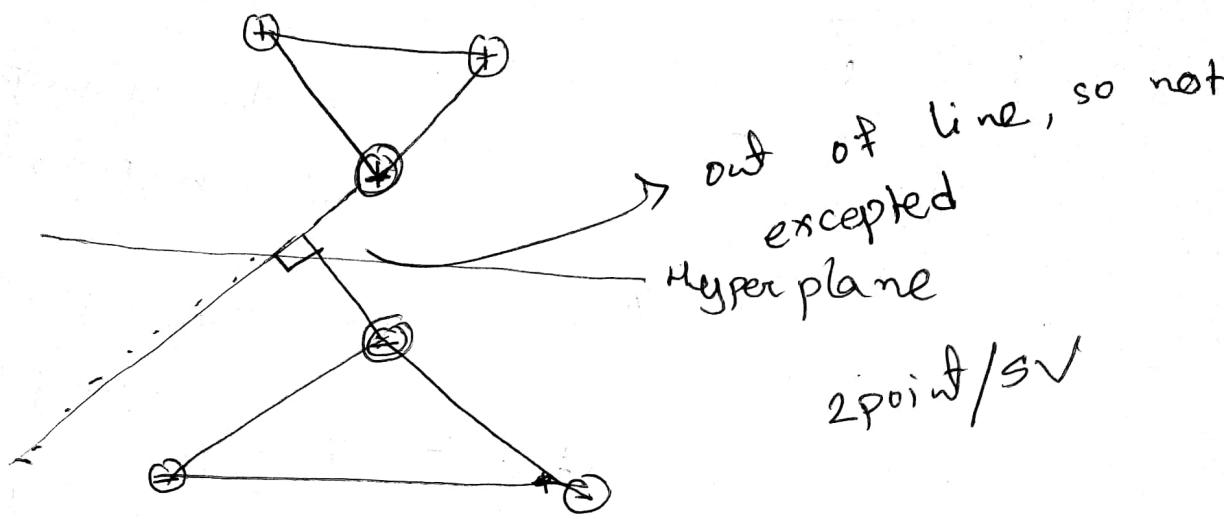
then + [b, c ,
constant]

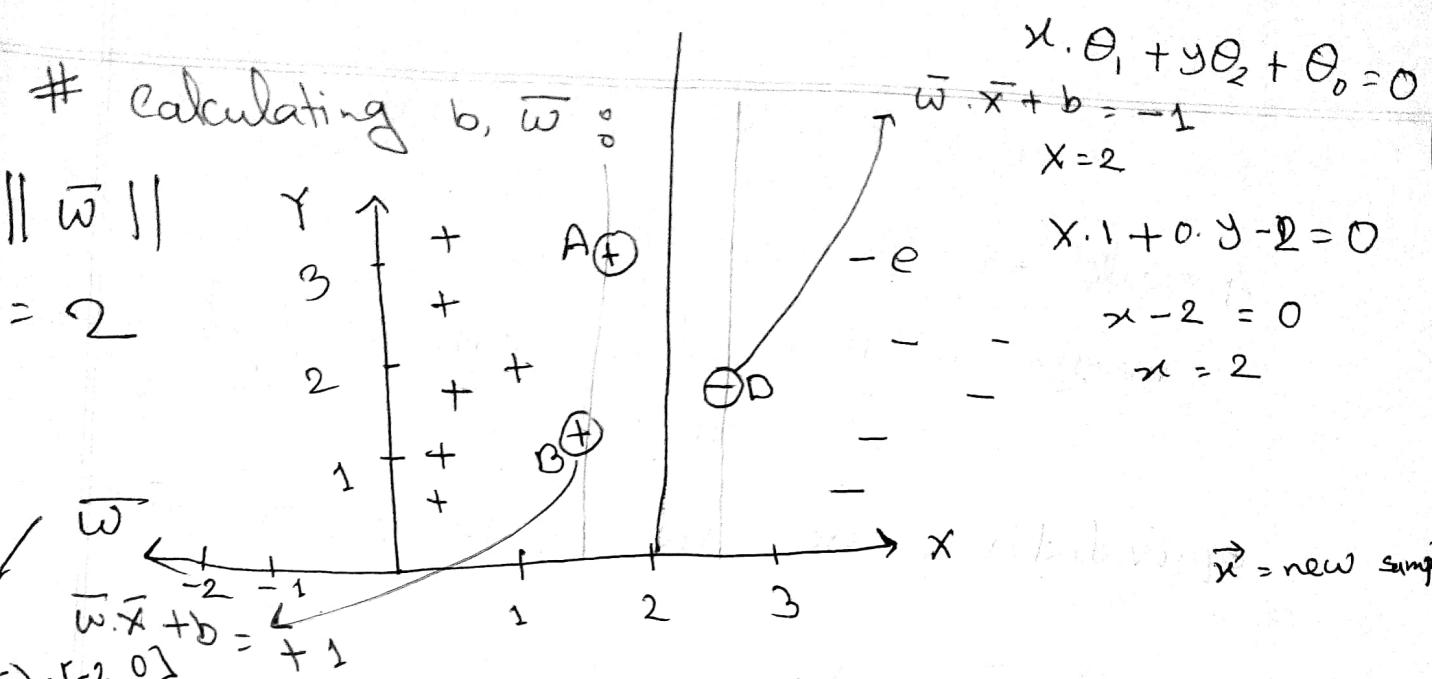
$$\text{or, } \bar{w} \cdot \bar{u} + b \geq 0$$

α = Supportiveness

Convex hull Method:







$w(\bar{w}) = [-2 \ 0]$

1. Draw the decision boundary. (DB)

2. Write the equation for the DB.

$$x = 2$$

$\begin{cases} 2 \text{ dimensional vector} \\ \text{as here two features} \end{cases}$

3. Rewrite the equation as $\bar{w} \cdot \vec{x} + b = 0$ — ①

$$\bar{w} \cdot \vec{x} + b = 0 \quad \left[\begin{array}{c} x \\ y \end{array} \right]$$

$$x \cdot \theta_1 + y \cdot \theta_2 + \theta_0 = 0$$

$$\theta^T x$$

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

$$\bar{w} \cdot \vec{x} + b = 0$$

$$\left[\begin{array}{c} 1 \\ 0 \end{array} \right] \left[\begin{array}{c} x \\ y \end{array} \right] + (-2) = 0$$

$$\text{class}(\vec{x}) = \text{sign}(\vec{\omega} \cdot \vec{x} + b)$$

4. Scale equation so that,

$$\text{margin width} = \frac{2}{\|\vec{\omega}\|} \text{ AND } \vec{\omega} \text{ points toward } + \quad \text{--- (2)}$$

equivalently, use the gutter constraint.

$$\vec{\omega} \cdot \vec{x}_i = y_i \quad / \quad \vec{\omega} \cdot \vec{v} + b = y_i \quad \text{--- (3)}$$

for all support vectors

where $y_i = \text{class}(i) = +1 \text{ or } -1$

$$C. \vec{\omega} \cdot \vec{x}_i + b = y_i$$

Using SV D,

$$\vec{\omega} \cdot \vec{x} + b = 0$$

$$C. [1 \ 0] \begin{bmatrix} 2.5 \\ 2.25 \end{bmatrix} + C(-2) = y_i = -1$$

$$2 - 2.5 = -1$$

$$-0.5 = -1$$

$$C \cdot 2 - C \cdot (1.5) = 1$$

$$\Rightarrow 0.5 C = -1$$

$$\Rightarrow C = -2$$

$$\left| \begin{array}{l} \vec{\omega} = \vec{\omega} \cdot C \\ = [-2 \ 0] \\ b = +3 \end{array} \right.$$

$$\text{margin}, \omega = \frac{2}{\|w\|} \quad \rightarrow \quad \sqrt{(-2)^2 + 0^2}$$

Euclidean dist. of 2

$$\text{margin width} = \frac{2}{\|w\|} = \frac{2}{2} = 1$$

$$= 1$$

→ if margin decreases,
w will increase

else, decrease

if margin decreases, w will increase

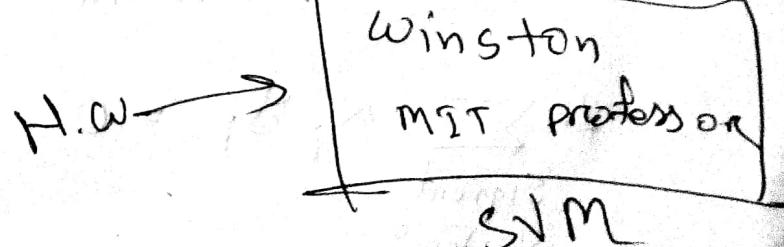
if margin increases, w will decrease

$$\text{margin width} = \frac{2}{\|w\|}$$

$$\text{margin width} = \frac{2}{\|w\|}$$

$$\text{margin width} = \frac{2}{\|w\|}$$

$$\text{margin width} = \frac{2}{\|w\|}$$

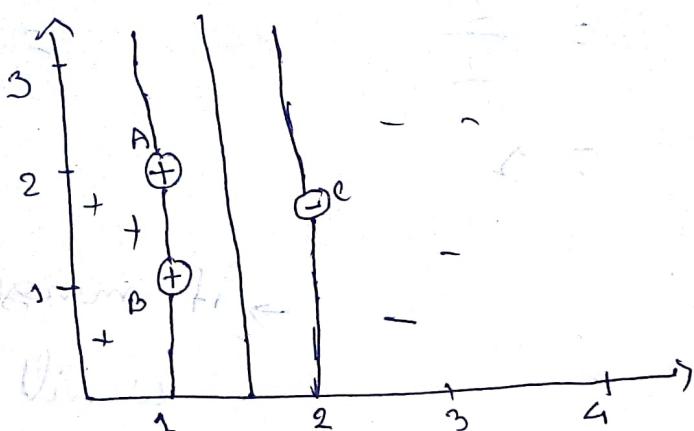


11.09.18

(Sequential minimal optimizers)
Calculating α (Lang Lagrange multipliers) /

Supportiveness value %

↓
find local min/max of
a function



* α represents how big of a role each point plays in determining / supporting the decision boundary.

1. For all non-SV support vectors, $\alpha = 0$

otherwise, the point is a SV and lies on the gutter and $\alpha > 0$

[Positive Supportiveness]

$$2. \sum_{\text{Support Vectors}} \alpha_i y_i = 0$$

Support
Vectors

$$\Leftrightarrow \sum_{\substack{\text{Positive} \\ S \vee P}} \omega_p = \sum_{\substack{\text{negative} \\ S \vee n}} \omega_n$$

$$\omega_A + \omega_B - \omega_C = 0 \Rightarrow \omega_A + \omega_B = \omega_C \quad \text{--- (1)}$$

$$3. \sum_i \omega_i y_i \bar{x}_i = \bar{\omega}$$

$$\Rightarrow \bar{\omega} = \sum_P \omega_P \bar{x}_P - \sum_n \omega_n \bar{x}_n$$

$$\begin{bmatrix} -2 \\ 0 \end{bmatrix} = \omega_A \begin{bmatrix} \frac{1}{3} \end{bmatrix} + \omega_B \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \omega_C \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

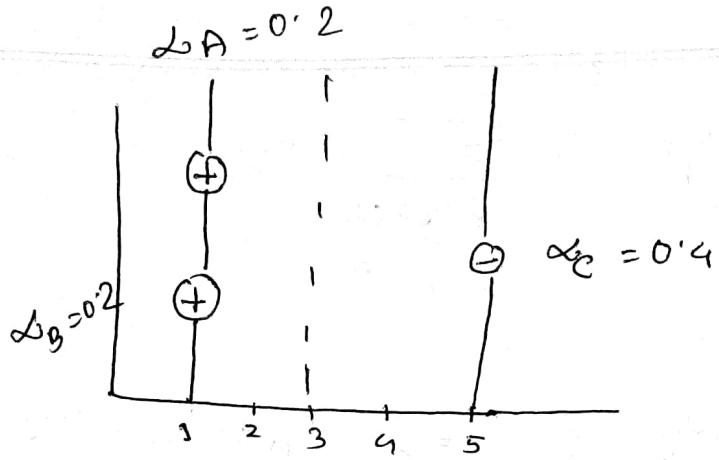
$$\Rightarrow -2 = \omega_A \cdot 1 + \omega_B \cdot 1 - \omega_C \cdot 2 \quad \text{--- (2)}$$

$$0 = \omega_A \cdot 3 + \omega_B \cdot 1 - \omega_C \cdot 2 \quad \text{--- (3)}$$

$$\omega_A = 1$$

$$\omega_B = 1$$

$$\omega_C = 2$$



Margin width = $\frac{2}{|w|}$

w ↓

Δ ↓

Margin width ↑

①

A (+)

Δ_C = 2

②

B (-)

Δ_B = 1.8

Δ min
width max

↓
Select

③

A (+)

④

A (+)

-c

B (+)

When data is not linearly separable,

- Apply transformation
- Replace dot product with another function (Kernel), effectively transforming the space.

(A Kernel function $K(\bar{u}, \bar{v})$ is a similarity measure)

$$\text{Class } (\bar{x}) = \text{Sign}(\bar{w} \cdot \bar{x} + b)$$

$$= \text{Sign}\left(\sum_i (\alpha_i y_i \bar{x}_i) \cdot \bar{x} + b\right)$$

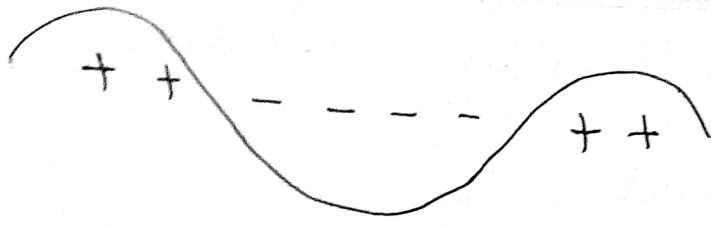
$$= \text{Sign}\left(\sum_i (\alpha_i y_i) K(\bar{x}_i, \bar{x}) + b\right)$$

↳ Transforming
to higher order/
function/dimension

Common kernel function:

- linear, e.g. $K(\bar{u} \cdot \bar{v}) = \bar{u} \cdot \bar{v} + 1$; draw line

- quadratic, e.g. $K(\bar{u} \cdot \bar{v}) = (\bar{u} \cdot \bar{v} + 1)^2$; draw conic Δ

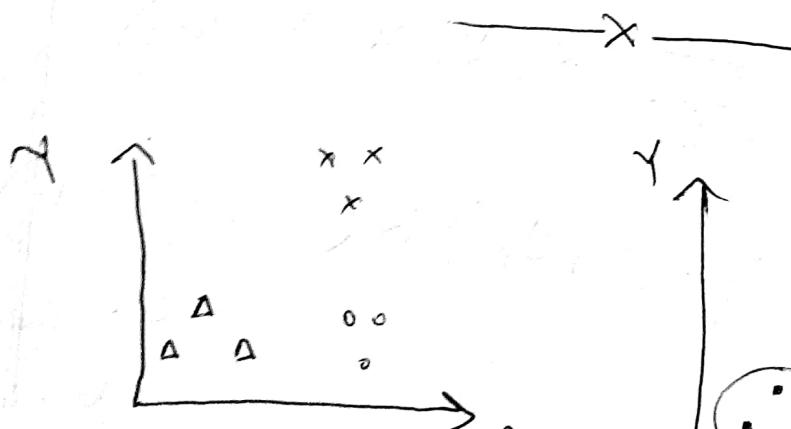


$$+ \oplus 0 \oplus + \oplus 0 0 \oplus 0 \oplus 0 \oplus 0 \oplus$$

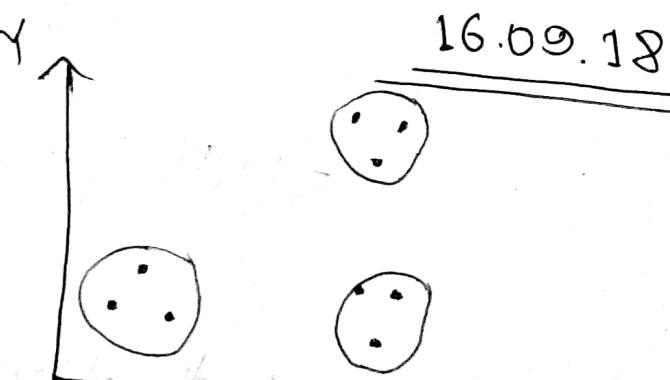
→ Polynomial

- Radial Basis Function Kernel

$$k(\bar{u}, \bar{v}) = e^{-|\bar{u} - \bar{v}|^2/6}$$



Supervised



Unsupervised (not
labeled)

⇒ Supervised learning:

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

⇒ Unsupervised learning:

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

Applications of clustering algorithms:

→ Social Network analysis

→ Market segmentation

→ Astronomical data analysis

* K-means Algorithm:

K = No. of clusters.



→ Select Random point k_1, k_2 on A, B.

* K-means:

→ iterative algo

→ has two steps

- (i) Cluster assignment step
(ii) Move centroid step.

repeat

K-means Algorithm

Input:

- K (# of cluster)

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

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

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

Dimension \mathbb{R}^{n+1} .

* Algorithm:

→ Randomly initialize K -cluster centroids

$$\mu_1, \mu_2, \dots, \mu_K$$

→ Repeat {

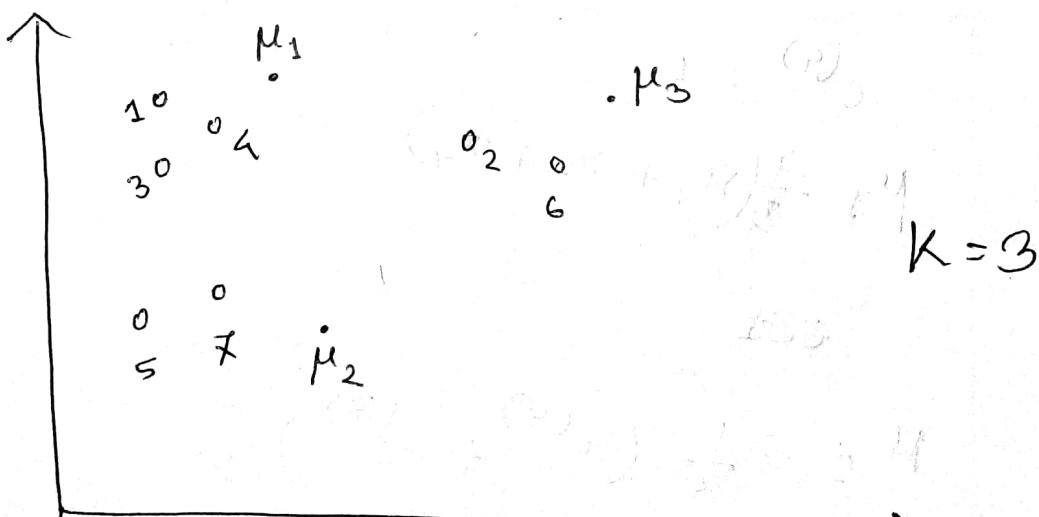
Clusters
Assignment
Step

min & $J(-)$ w.r.t.
 $c^{(1)}, c^{(2)}, \dots, c^{(m)}$
 $\mu_1, \mu_2, \dots, \mu_K \rightarrow$ fixed

for $i = 1$ to m
 $c^{(i)} =$ index (1 to K) of cluster centroid
closet to $(x^{(i)})$

for $k = 1$ to K
 $\mu_k =$ average (mean) of points assigned
to cluster K .

move
centroid
step



$$K=3$$

$$x^{(1)} \Rightarrow c^{(1)} = 1$$

$$x^{(2)} \Rightarrow c^{(2)} = 3$$

$$x^{(3)} \Rightarrow c^{(3)} = 1$$

$$x^{(4)} \Rightarrow c^{(4)} = 1$$

$$x^{(5)} \Rightarrow c^{(5)} = 2$$

$$x^{(6)} \Rightarrow c^{(6)} = 3$$

$$x^{(7)} \Rightarrow c^{(7)} = 2$$

$$c^{(1)} = 1$$

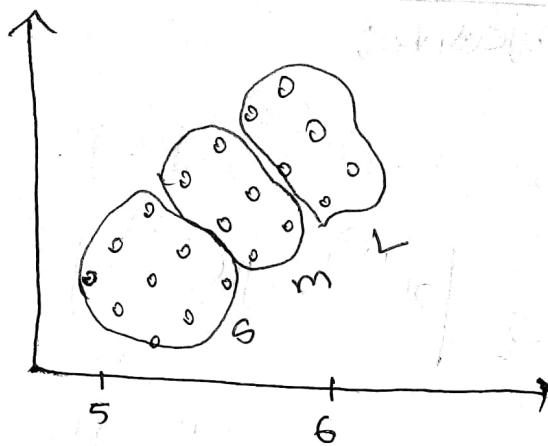
$$c^{(3)} = 1$$

$$c^{(9)} = 1$$

$$\mu_1 = \frac{1}{3}(x_1 + x_3 + x_9)$$

~~200~~

$$\mu_2 = \frac{1}{2} (x^{(5)} + x^{(7)})$$



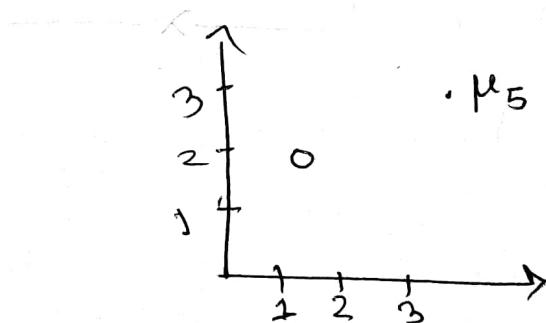
if T-shirt size

Optimization objective:

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

μ_k = cluster centroid $K (\mu_k \in \mathbb{R}^n)$
(same with the sample dimension)

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



$$x^{(1)} = (1, 2)$$

$$\mu_5 = (2, 3)$$

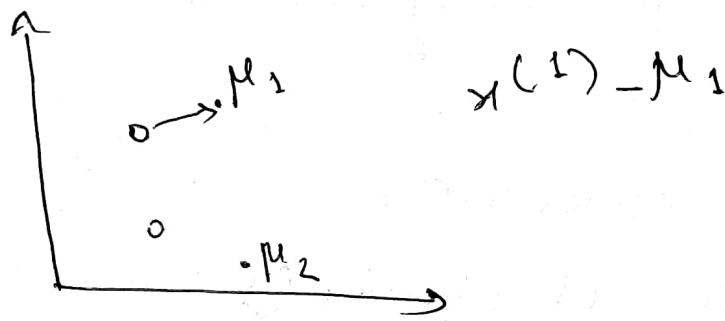
$$\therefore c^{(1)} = 5$$

$$\mu_c^{(1)} = \mu_5, (2, 3)$$

* Optimization objective:

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

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



$$\min c^{(1)}, c^{(2)}, \dots, c^{(m)}$$

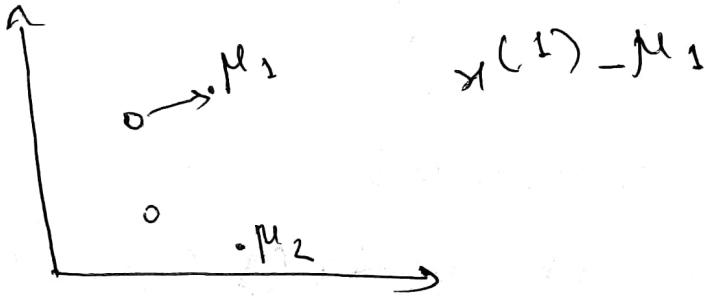
$$\mu_1, \mu_2, \dots, \mu_k \quad J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_k)$$

\underline{x}

* Optimization objective:

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

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



$$\min c^{(1)}, c^{(2)}, \dots, c^{(m)}$$

$$\mu_1, \mu_2, \dots, \mu_k$$

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

— x —

18.09.18

K-mean Algo.:

- Randomly initialize k -clusters centroid

$$\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$$

- Repeat {

for $i = 1$ to m

$c^{(i)}$ = index of cluster centroid (1 to k)
which is closest to $x^{(i)}$

for $K = 1$ to K

μ_K = average (mean) of points
assigned to cluster μ_K

Optimization:

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

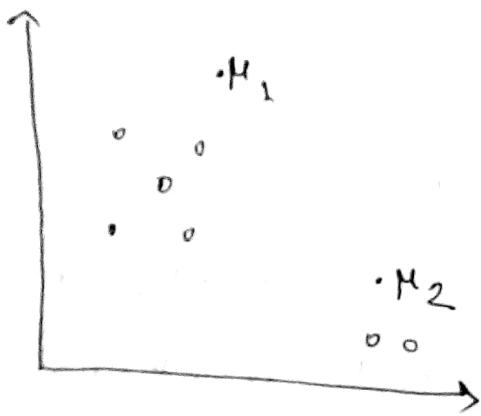
$$\min c^{(1)}, c^{(2)}, \dots, c^{(m)} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\mu_1, \mu_2, \dots, \mu_k$$

i^{th} index point $\rightarrow \mu_5$

$$c^{(i)} = 5$$

$\mu_{c^{(i)}} = \mu_5 = \text{Position value}$



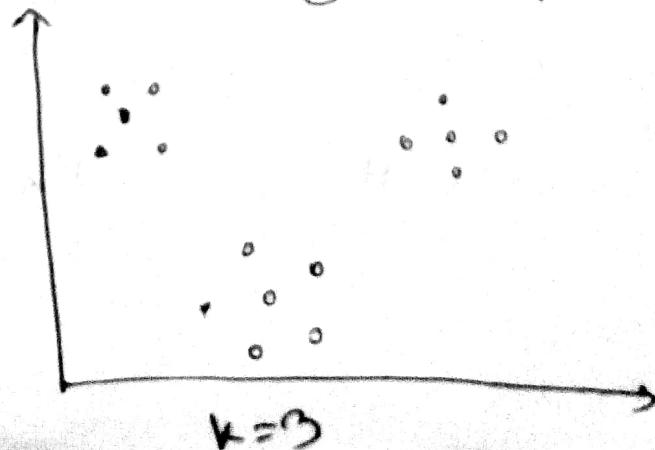
Important:

* assignment of each point

* Position selection of μ_k

Random Initialization:

- Select k such that $k \leq m$.
- Randomly pick k training examples.



- Set $\mu_1, \mu_2, \dots, \mu_K$ equal to these examples.

Drawbacks:

* Local optimal:

Random initialization to avoid local optima

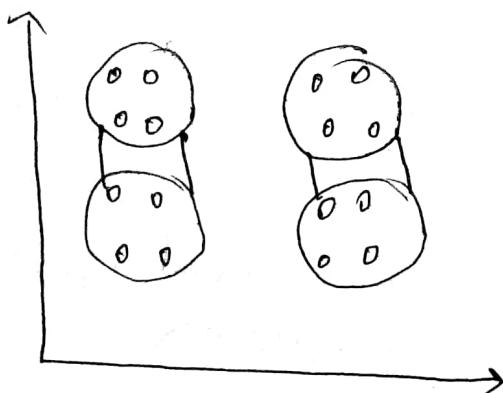
- for $i = 1$ to 100 ($50 - 1000$)
 - { • Randomly initialize K cluster centroid.
 - Run K -means algorithm
 - Get $c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_K$
 - Compute cost function (distortion)
- $$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$
- }

Pick the cluster that gives the lower cost value.

* If $K = 2$ to 10 this algo. works really well.

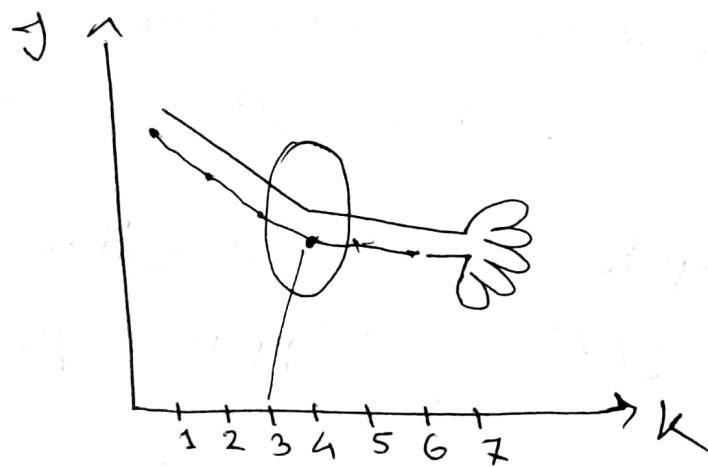
* If $K = 100$ + expected result can be get within few iterations. So, no. of iterations can be decreased.

How to choose value of K ?

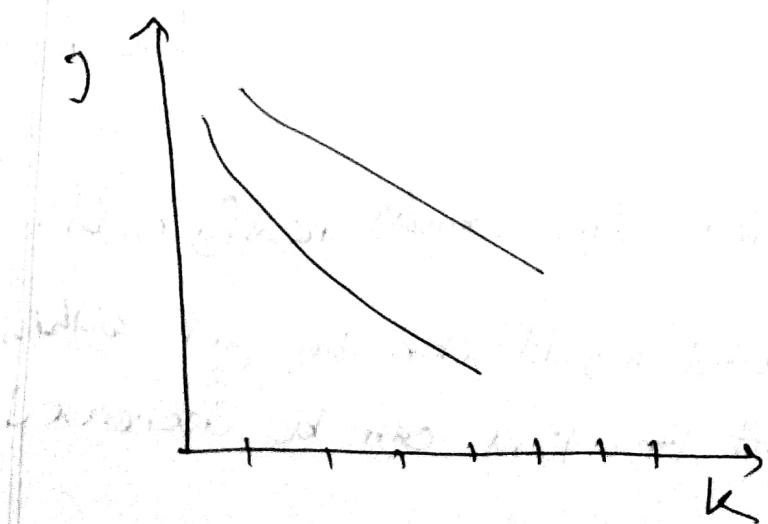


$$k = 4 / 2$$

* Elbow Method :



↑ Select 3 as ~~the~~ the value of k.



→ Elbow method

doesn't work
well in these
kinds of cases

* choosing K's value based on purpose

