

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

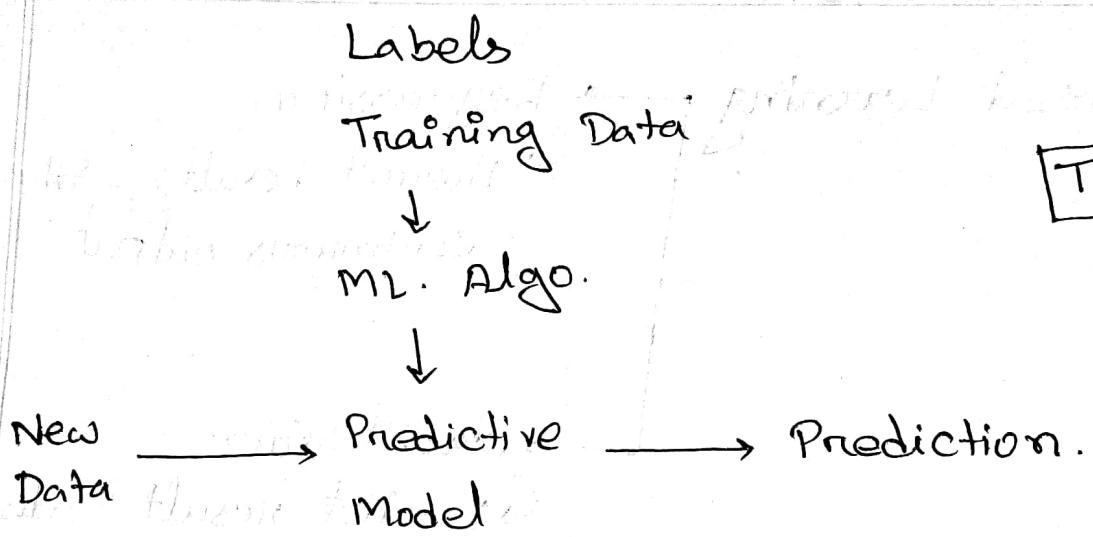
Acenos Vs. Non-Acenos

1. Supervised Learning:

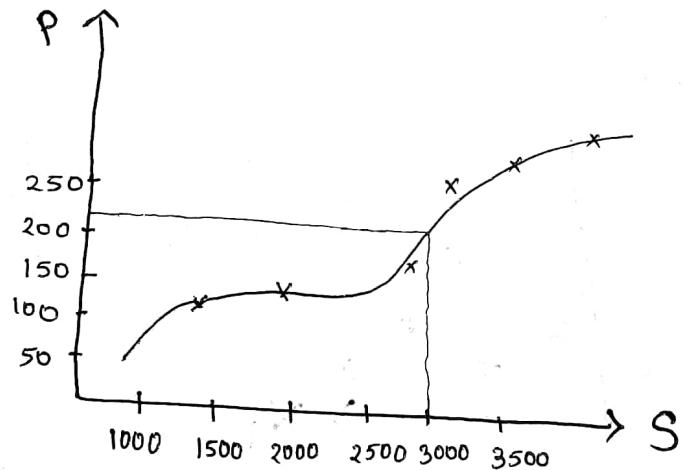
SL is where you have input variables (x) and output variables (y) and you use an algorithm to learn the mapping function from the input to the output.

$$Y = f(x) \text{ for } (x, y) \text{ (label data)}$$

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.

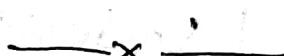


Model



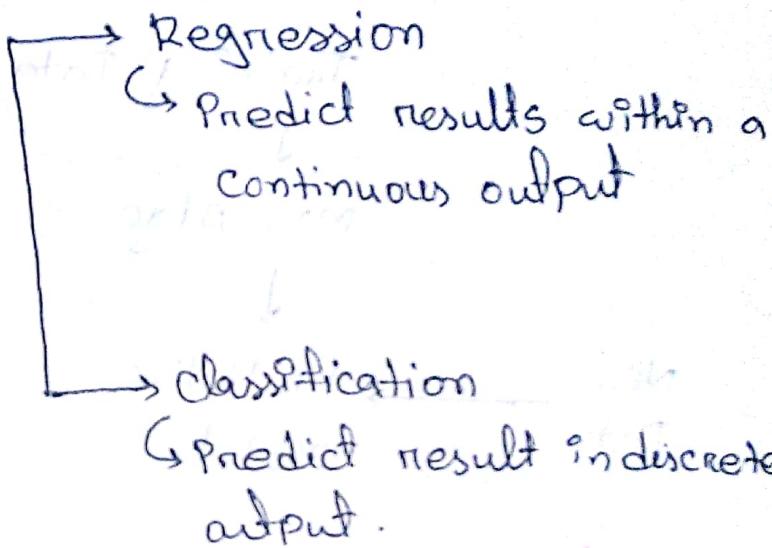
$$y = f(x)$$

$$P = f(S)$$

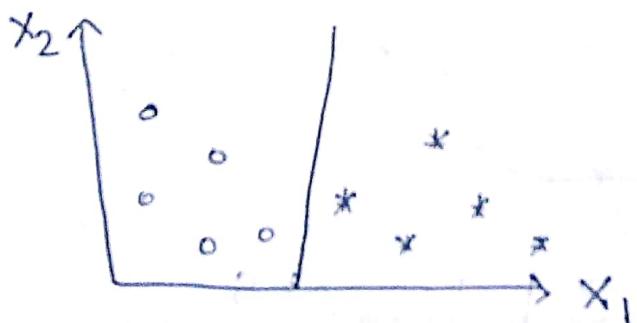


08.05.18

Supervised Learning



* Classification:



* Unsupervised Learning:

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

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

Unsupervised Learning

(algo) makes no explicit
as based organization
and communities self

Clustering

↳ want to discover the
inherent grouping in the
data such as group
customers by purchasing
behaviors. (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 teacher yet.

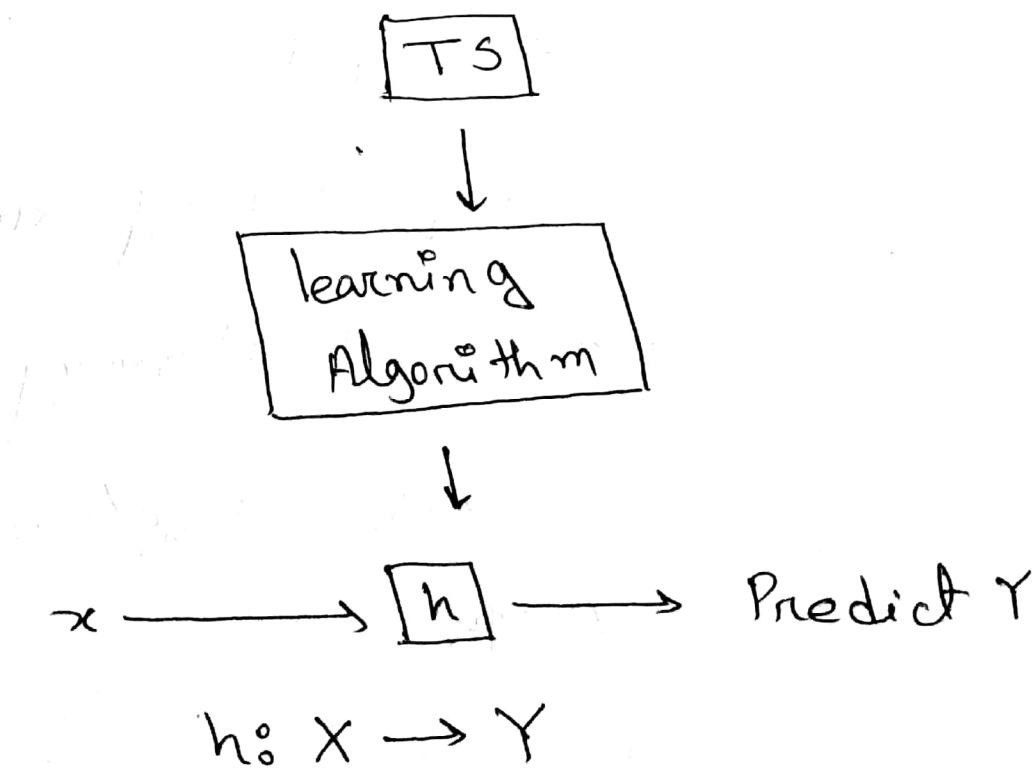


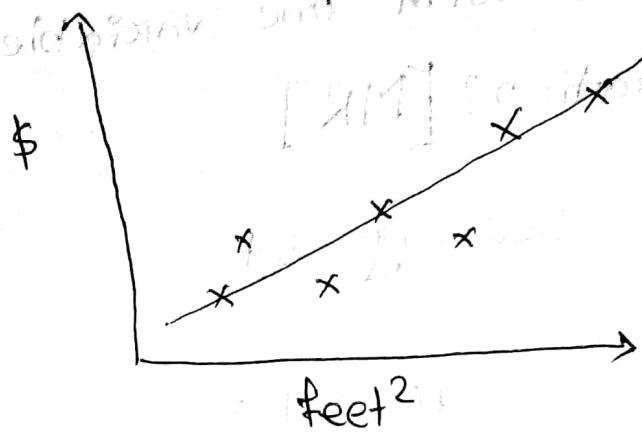
* linear regression with one variable:

Model Representation: [MR]

TS → Training Set

<u>Feet²</u>	<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

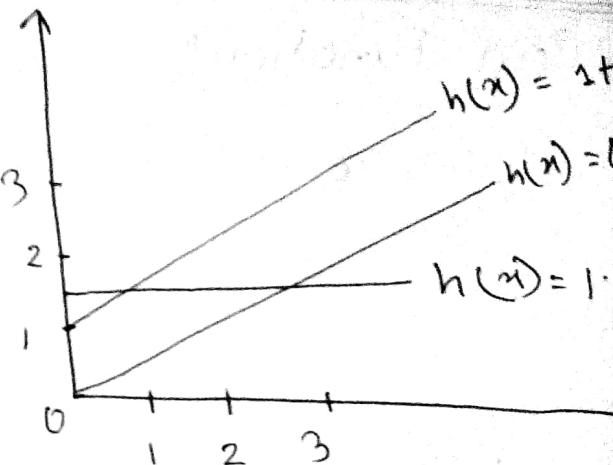
$(x^{(i)}, y^{(i)})$ = i^{th} training example

$\theta \rightarrow$ Parameters/weights

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

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

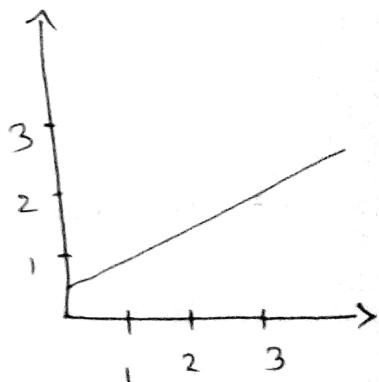
$$\begin{array}{c|c|c} \theta_0 = 1.5 & \theta_0 = 0 & 1 \\ \hline \theta_1 = 0 & \theta_1 = 0.5 & 0.5 \end{array}$$



\checkmark

$$\begin{array}{ccccc} \theta_0 = & 0 & 0.5 & 1 & 1 \\ \theta_1 = & 1 & 1 & 0.5 & 1 \end{array}$$

<u>i.e.</u>	X	Y
	0	4
	1	7
	2	7
	3	8



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

$$= 2 + 2 * x$$

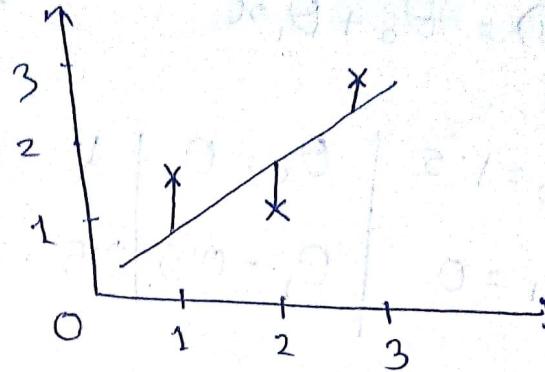
$$\text{if } x = 0 \rightarrow 2 + 2 \cdot 0 = 2$$

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

* Cost Function:

minimize θ_0, θ_1

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

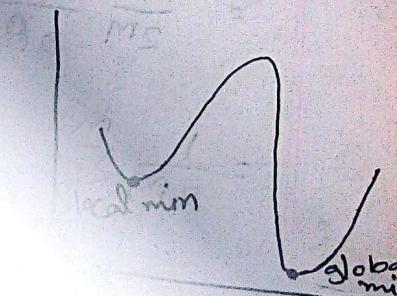
Cost function:

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

Goal: minimize $\theta_0, \theta_1 \quad J(\theta_0, \theta_1)$

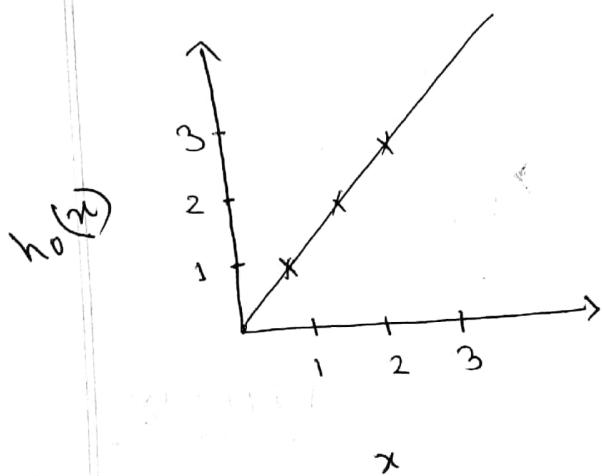
(any value)

until we have min



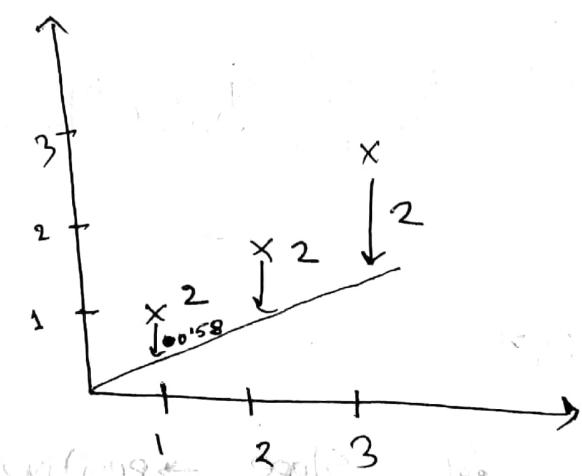
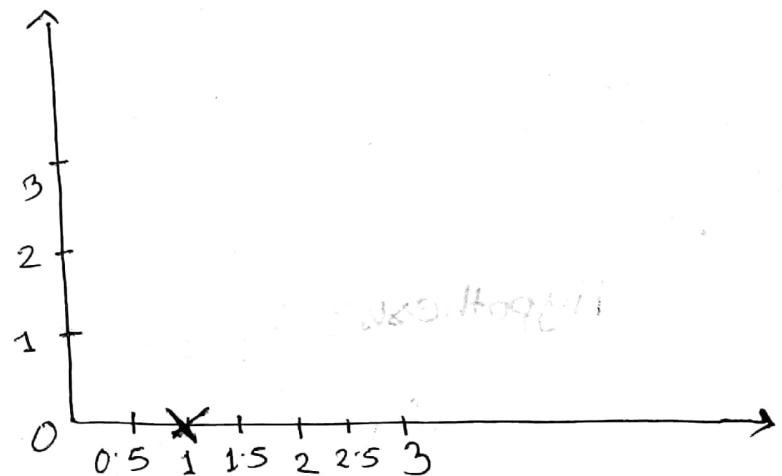
$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \text{let, } \theta_0 = 0$$

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



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

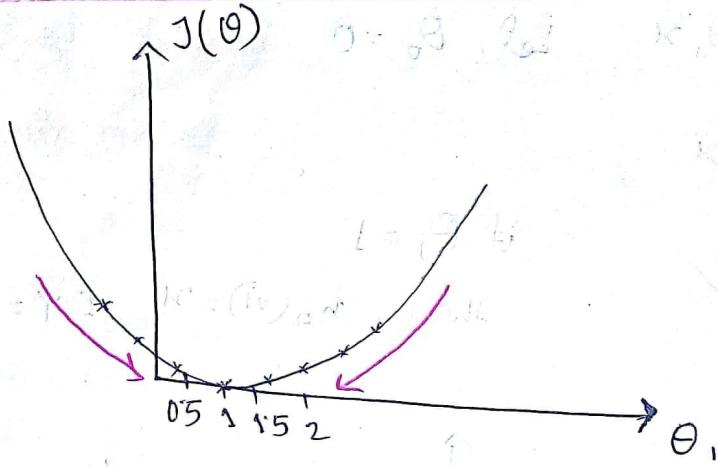
$$\text{then, } h_{\theta}(x) = x \quad [Y = x]$$



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

$$\theta_1 = 0.5$$

$$\begin{aligned} J(0.5) &= \frac{1}{2m} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \\ &= \frac{1}{2 \times 3} () = 0.58 \end{aligned}$$



15.05.18

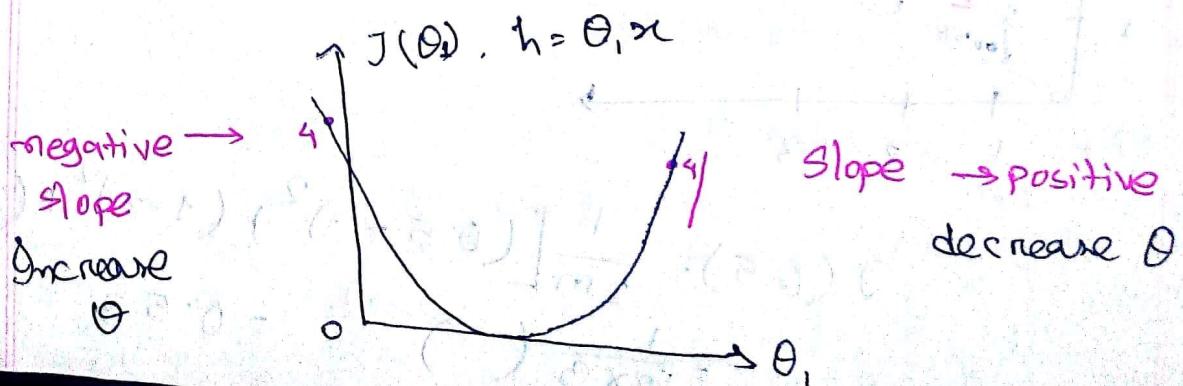
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

$\theta_1, \theta_0 \rightarrow$ 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



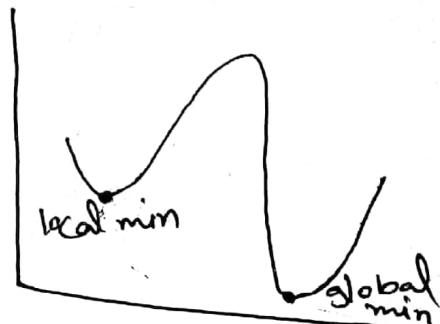
Gradient Descent:

- 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

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

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

when $j = 0$,

$$\begin{aligned}\theta_0 &\rightarrow \frac{\partial}{\partial \theta_0} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\&= \frac{1}{2m} \frac{\partial}{\partial \theta_0} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\&= \frac{1}{2m} \sum_{i=1}^m 2 \cdot (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_0} (h_\theta(x^{(i)}) - y^{(i)}) \\&= \frac{1}{2m} \cdot 2 \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_0} (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \\&= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}). 1\end{aligned}$$

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

when, $j = 1$,

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

$$\begin{aligned}
 &= \frac{1}{2m^2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_1} (h_\theta(x^{(i)}) - y^{(i)}) \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_1} (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}
 \end{aligned}$$

repeat

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

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

update simultaneously:

$$\left\{ \text{temp} \theta = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta)
 \right.$$

$$\text{temp}^1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta)$$

$$\theta_0 = \text{temp} \theta$$

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

}

Batch Gradient descent:

Solution

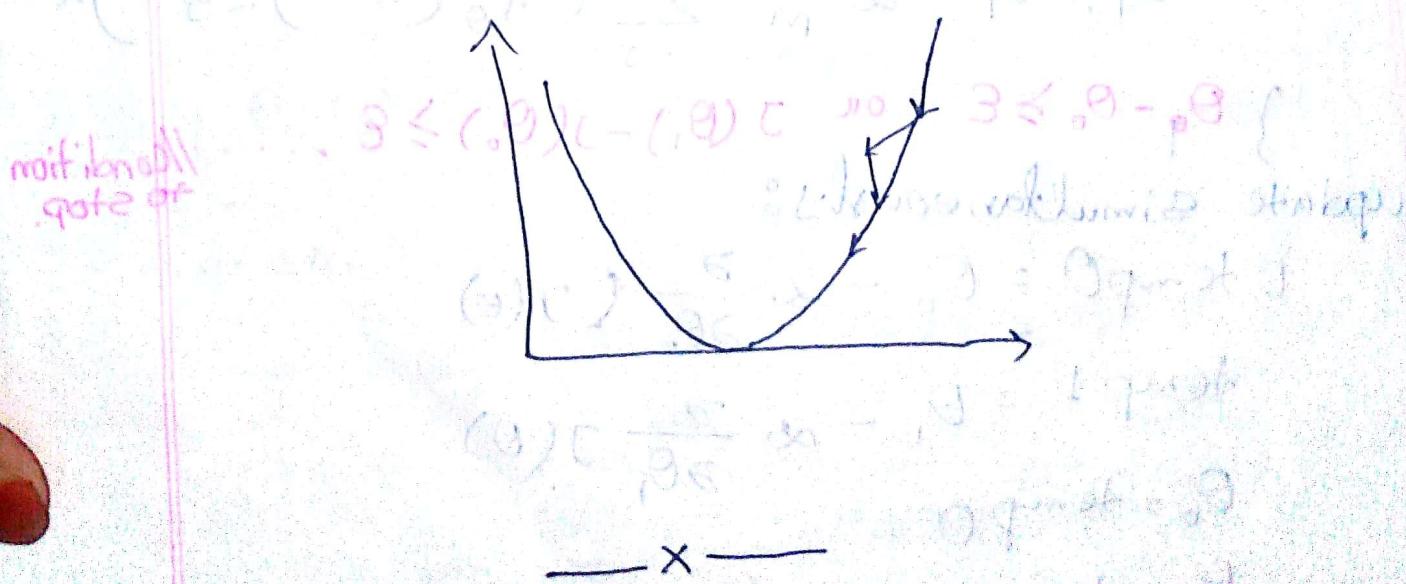
Incremental / Stochastic Gradient descent:

(Fast, but possibility of getting global solution is very low.)

Loop {

for i = 1 to m {

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



"Special class"

17.05.18.

eliademy.com

~~Python~~

~~mr. jordan~~

Introduction to Machine Learning and Decision Tree

aktarhossain@daftodilvarity.edu.bd

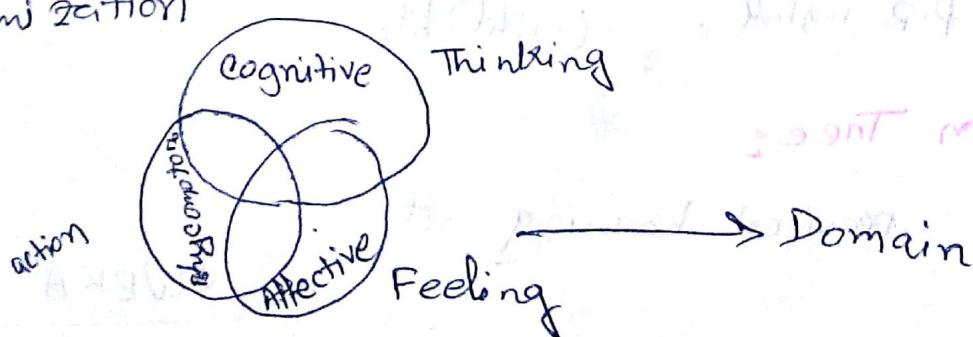
Industry 4.0

* Three components of ML

Representation

Evaluation

Optimization



Bloom's Taxonomy of learning

Supervised or inductive learning:

markany.com

animal objects → Power
2**100 most abundant

sentences
every single type is tuple

Python Basics.pdf.

Everything in Python is an object.

Object

1. get-pip.py

Install pip installer

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

Palindrome

2. Install Numpy & sci-pi

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

Decision Tree:

Supervised learning set

"WEKA" → DM Tool

New Zealand Waikato University

	x_0	x_1 feet ²	x_2	x_3	x_4	20.05.18.
		# bed rooms		# floors	Age	Price
1	2104	5		1	45	460
1	1416	3		2	40	232
1	1535	3		2	30	315
1.	852	2		1	36	178

$n = \# \text{ features}$

$$x^{(i)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \rightarrow \begin{array}{l} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$$

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

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

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

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

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

Let, $x_0 = 1$ for all $i=1 \dots m$

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

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

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

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

$n+1$ element 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:

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

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

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

when, $n = 1$ *gradient descent* - the simple #

Repeat {

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

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

} *Simple gradient descent*

when $n \geq 2$

Repeat {

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

}

$$\theta' \geq (\theta_x) \geq \theta' -$$

Update Simultaneously

Assignment
Derivation of
gradient descent
next class

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

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

$$\vdots$$

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

Gradient Descent in practice

→ Feature Scaling

→ Mean normalization

→ Debug

→ select α

$$x_1 = \text{size} = 100 - 2500 \text{ feet}^2$$

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

* θ will descent quickly on small range
So, set input value in same range (roughly),

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

$$\begin{aligned} -0.5 &\leq (x^{(i)}) \leq 0.5 \\ \boxed{-3 \leq x^{(i)} \leq 3} \end{aligned}$$

Feature Scaling:

$$x_j = \frac{x_j}{\text{range}} = \frac{2104}{2400} \approx 1 \quad \text{for } x_1$$

$$x_2 = \frac{5}{3} = 1.5$$

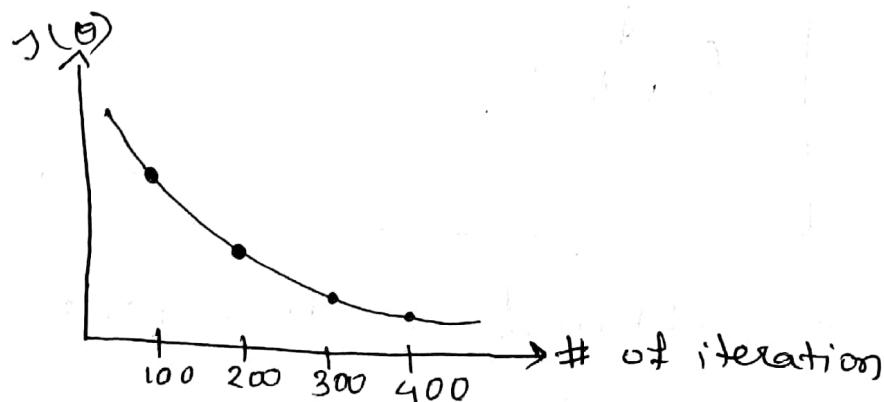
Mean normalization:

$$x_j = \frac{x_j - \mu_j (\text{mean avg.})}{s_j (\text{range})}$$

$$x_2^{(1)} = \frac{5 - 3.25}{3}$$

$$\mu_j = \frac{5 + 3 + 3 + 2}{4}$$

Debug:



$$J(\theta) = 5$$

2

$$J(\theta) = 4.99999$$

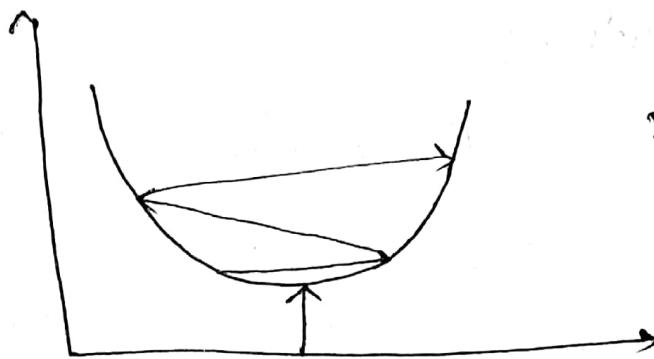
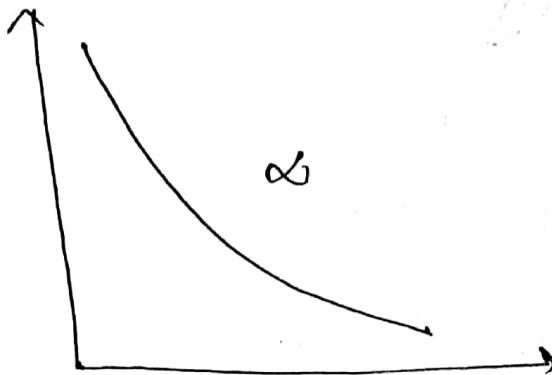
3

Change learning rate, α

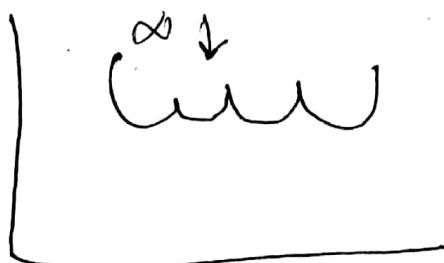
→ Threshold value $\approx 10^{-3}$

⇒ If the change is below threshold value then we can assume to get optimal solution.

Select α :



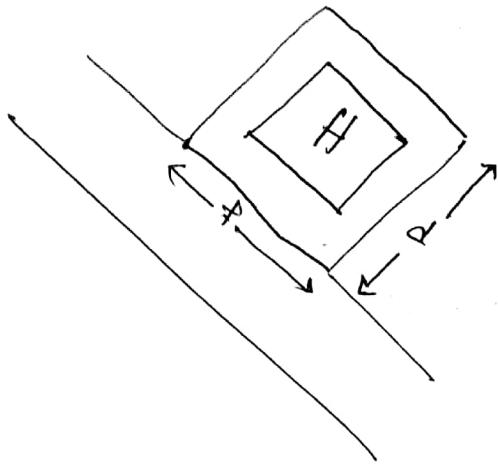
reduce the value
of α .



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

* 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
 \uparrow
 \uparrow

d
 f
 p

$$= \theta_0 + \theta_1 x$$

\uparrow
 \uparrow
 \uparrow

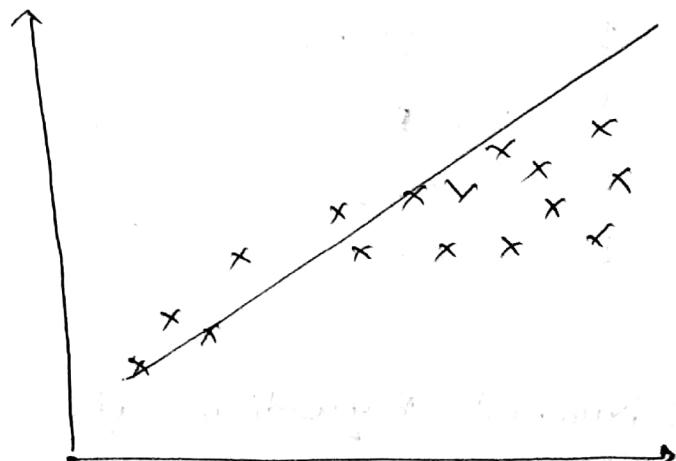
area

*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.

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

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

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

X = input matrix [showing four data instances]

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 90 \\ 1 & 1535 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 35 \end{bmatrix}$$

$$x^{(1)} = \begin{bmatrix} 1 \\ 2104 \\ 5 \\ 1 \\ 45 \end{bmatrix} \quad x^{(2)} = \begin{bmatrix} 1 \\ 1416 \\ 3 \\ 2 \\ 90 \end{bmatrix}$$

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

$$X \quad X^T$$

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

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. $O(n^3) \ll (X^T X)^{-1}$

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

4. Slow if n is large

[$n = \text{no. of features}$]

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

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

not invertible

On 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

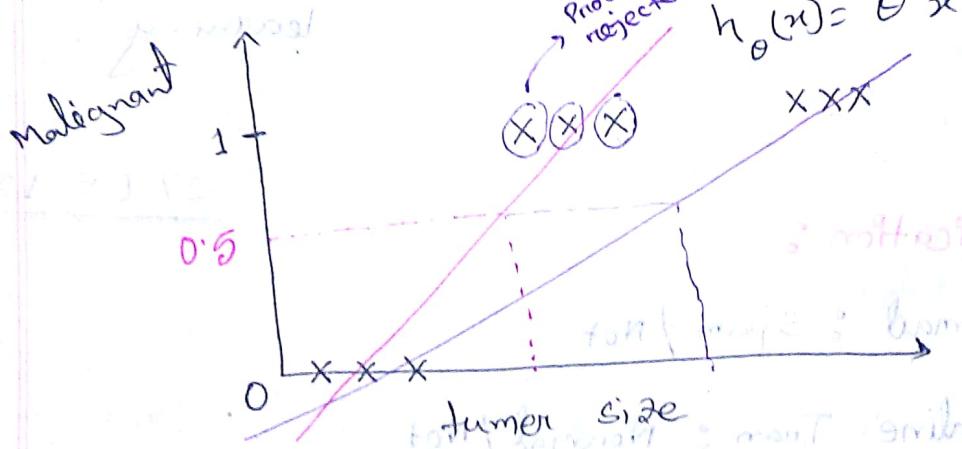
Something

1 : Positive class / Presence of something

* Problem with linear Regression for classification

Problem:

is a tumor malignant or not?



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

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

Problem:

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

$h_{\theta}(x)$ using linear regression may produce value

$$< 0 \text{ or } > 1$$

* Solution: \rightarrow

Logistic Regression (Logistic function, sigmoid function)
[classification algorithm]

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

$$\theta = (\theta_0, \theta_1)$$

$\frac{1}{1 + e^{-x}}$ as $x \rightarrow \infty$ then classification completed

always produces result -

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

Hypothesis:

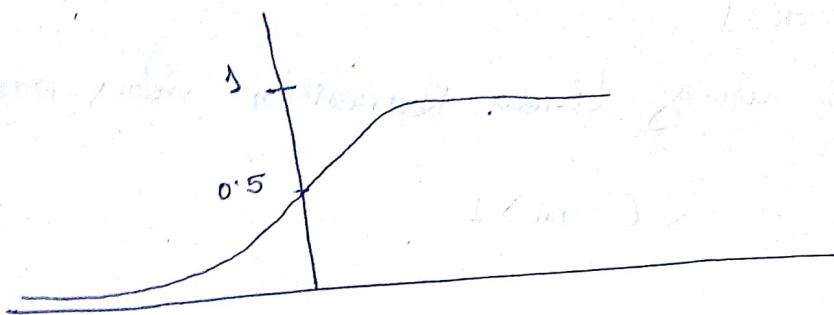
$$(0, x_1, x_2, \dots, x_n) \rightarrow h_{\theta}(x) = \theta^T x$$

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

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

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

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



$$g(2)$$

$$g(0) = 0.5$$

$1 - e^{-x}$ is positive &

$$\text{[and also } g(x) \text{ is odd]} \\ g(-\infty) = 0$$

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

\rightarrow Estimate Probability that $y=1$ on input x .

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

$h(x) = P(y=1 | x, \theta)$ Probability that $y=1$,

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

given x parameterized by θ

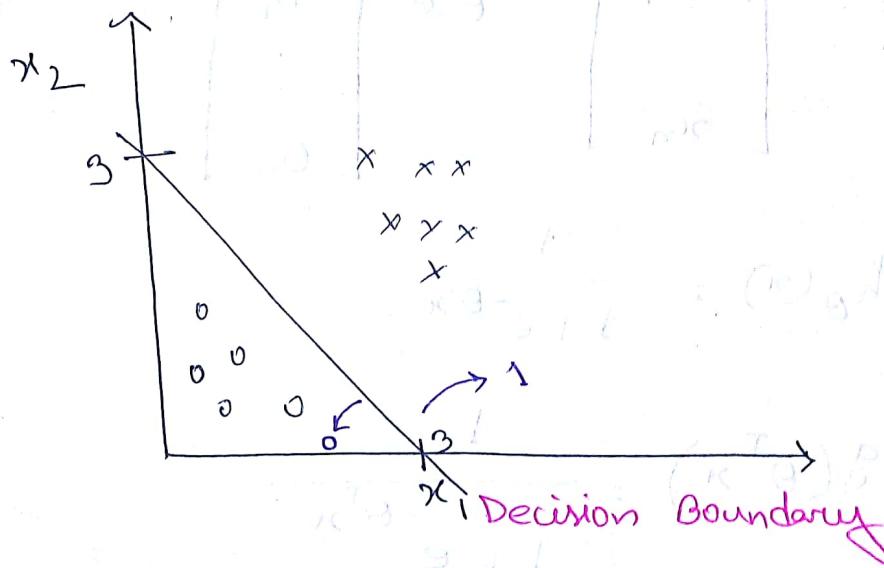
$$P(y=1 | x, \theta) + P(y=0 | x, \theta) = 1$$

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

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

Predict, $y = 1 \rightarrow h_0(x) \geq 0.5 \rightarrow \theta^T x \geq 0$
 $y = 0 \rightarrow h_0(x) < 0.5 \rightarrow \theta^T x < 0$

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



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

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

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

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

$$x_1 + x_2 = 3$$

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

29.05.18.

Logistic Regression: $1 = S, 0 = B$ Training Examples: $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ # Example: m training samples# Features: n number of features

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

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

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

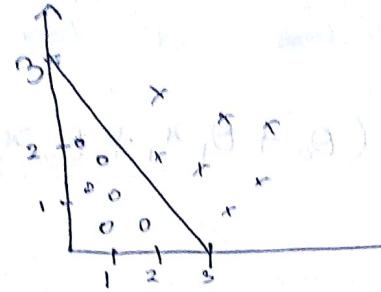
$$Y = 1, \quad g(2) \geq 0.5$$

$$\theta^T x = 0$$

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

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

2



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

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

$$\begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} = -3 + x_1 + x_2$$

if not satis. $\quad 1 = -3 + x_1 + x_2$

$y = 1 \text{ when } \theta^T x \geq 0$
 $y = 0 \text{ when } \theta^T x < 0$

constraint $x_1 + x_2 \geq 0$ is constraint satisfied
 $x_1 + x_2 \geq 0$ for given x to form

$x_1 + x_2 = 3$ to consider original #

given $x_1 + x_2 \geq 0$

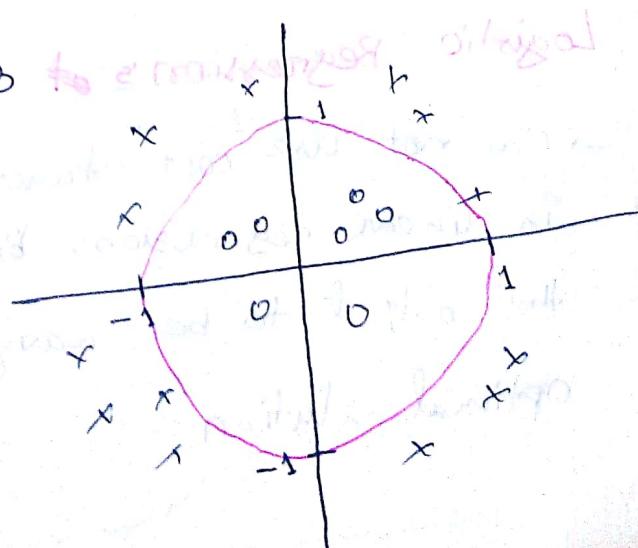
given $x_1 + x_2 = 3$

given $x_1 + x_2 \geq 0$

given $x_1 + x_2 = 3$

given $x_1 + x_2 \geq 0$

given $x_1 + x_2 = 3$



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 ~~not~~ 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 convex function.

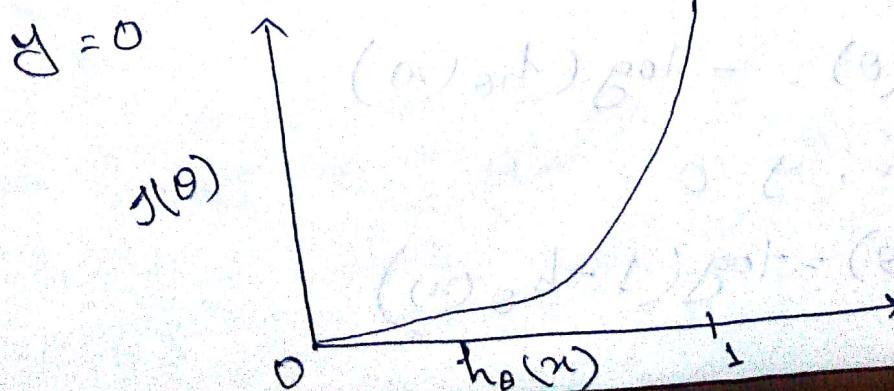
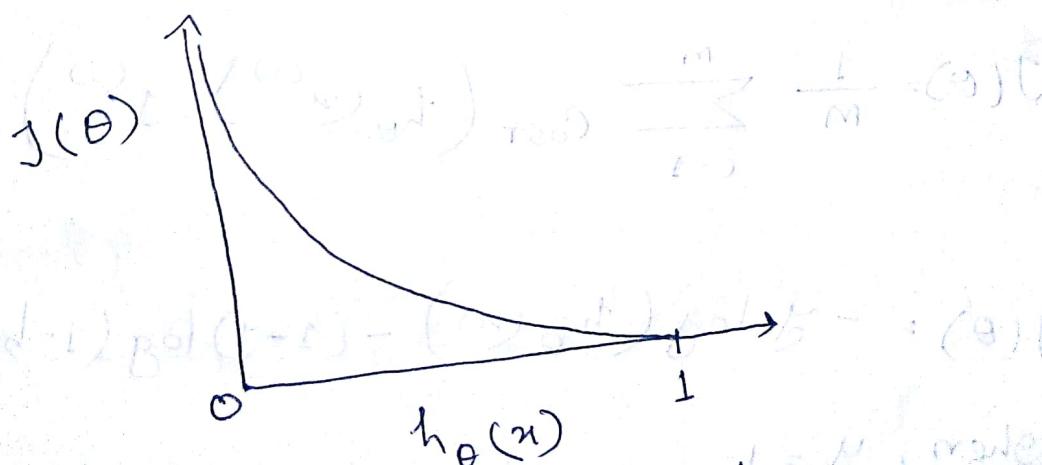
* Cost function:

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

$y = 1$ or 0

when

$$\hat{y} = 1 \quad \text{if } h_\theta(x) \geq 0.5 \quad \text{else } h_\theta(x)$$



5

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

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

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

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

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

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

when, $y=1$

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

when, $y=0$

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

6

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

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

Gradient Descent:

Repeat {
 $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$
 } // update simultaneously

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

Vector form: $\theta := \theta - \frac{\alpha}{m} \sum x^T (\phi(x\theta) - y)$

- BFGS

- L-BFGS

7

03.06.18.

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

$$g(z)$$

0.5

 $y = 1$

$$(d - o) \cdot \theta^T x \geq 0 \Rightarrow f(d, o, \theta) = 0$$

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

if true else if false

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

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

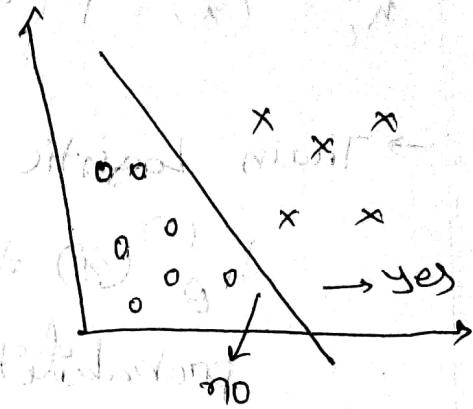
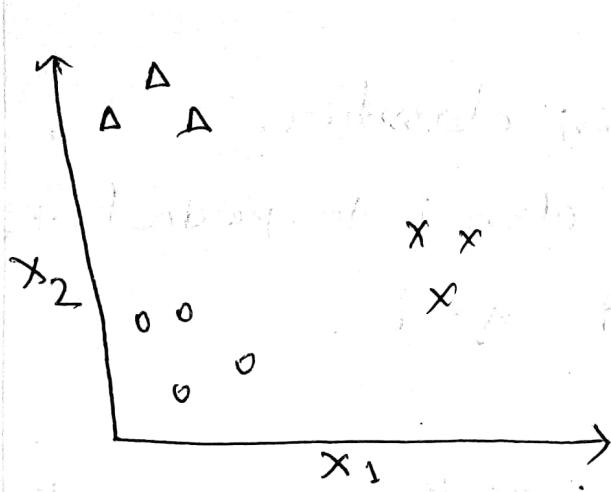
$y = h_{\theta}(x) ; \text{cost} = \text{zero}$

$$(f + (\theta^T x) \theta)^T x = \theta^T x^T x + \theta^T x = \theta^T x$$

cannot be zero

8

Multiclass classification:



multiclass classification

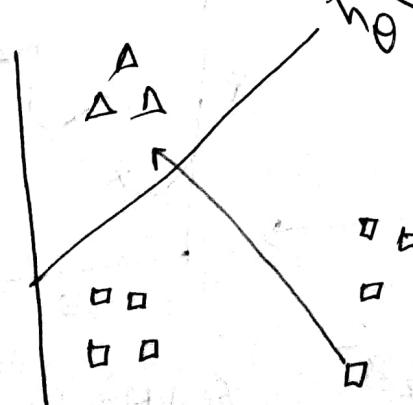
- * One vs all (one vs rest)

class 1: Δ

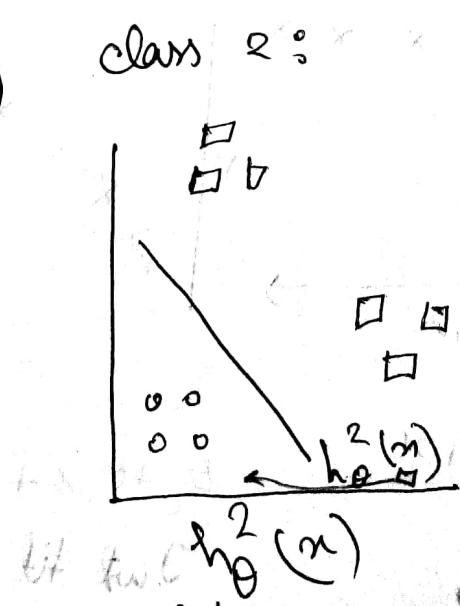
class 2: O

class 3: X

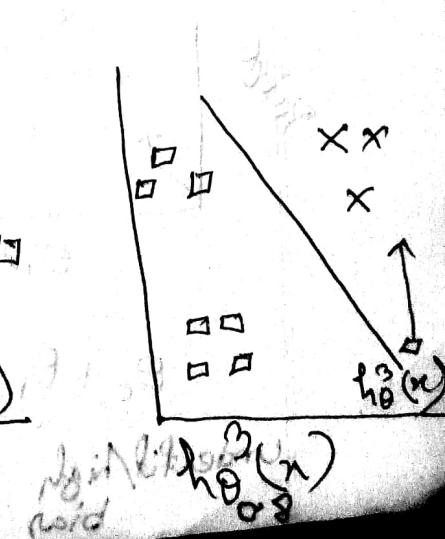
class 1: Δ $h_0^1(x)$



class 2: O $h_0^2(x)$



class 3: X $h_0^3(x)$



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

→ Train Logistic Regression classifier.

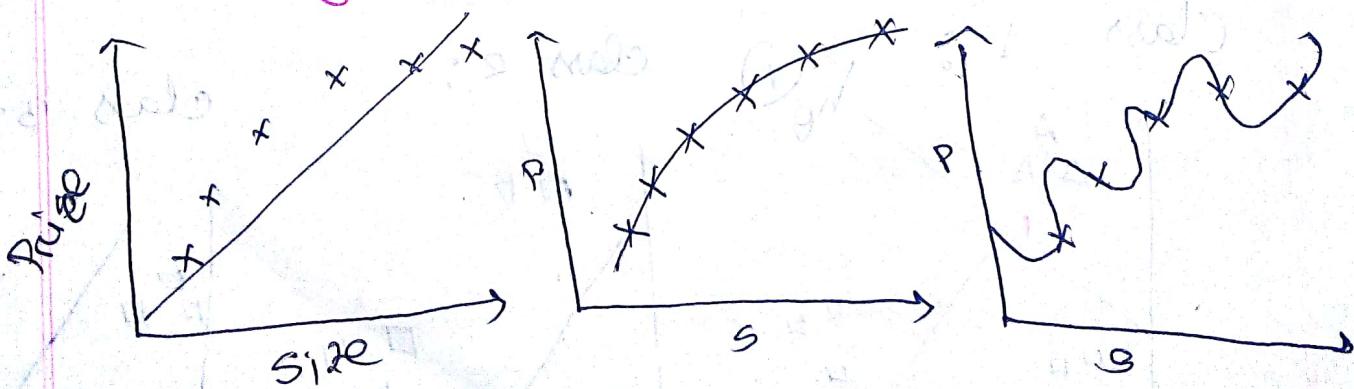
$h_{\theta}^{(i)}(x)$ for each class i to predict the probability that $y=i$.

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

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

— x —

Overfitting Vs Underfitting:



$$\theta_0 + \theta_1 x$$

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 problem:

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.

11

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

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

tends to $\theta \rightarrow 0$

0

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

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

λ = regularization constant.

x

R

05.06.18.

Regularization: (for linear Regression)

⇒ small value for θ s

→ simple hypothesis

→ less prone to overfitting

| if λ high,
 θ will be
small.

Features: $x_0, x_1, x_2, \dots, x_n$

P: $\theta_0, \theta_1, \theta_2, \dots, \theta_n$ not too large #
if λ is small, θ will be large

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

Gradient descent:

Repeat {

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

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

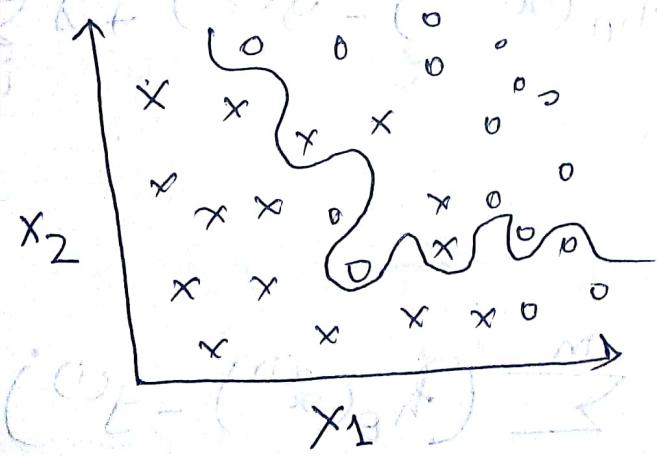
$$+ \frac{\lambda}{m} \theta_j \right]; j = 1, 2, 3, \dots, n$$

}

13

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

Regularization for logistic regressions



$$\begin{aligned}
 h_{\theta}(x) &= g(\theta^T x) \\
 &= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 \\
 &\quad + \theta_5 x_1 x_2 + \theta_6 x_1^2 x_2 + \dots)
 \end{aligned}$$

14

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

Gradient descent:

repeat {

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

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

}

 $\rightarrow x$

Function linear-Reg(theta-0, theta-1, x)

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

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

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)

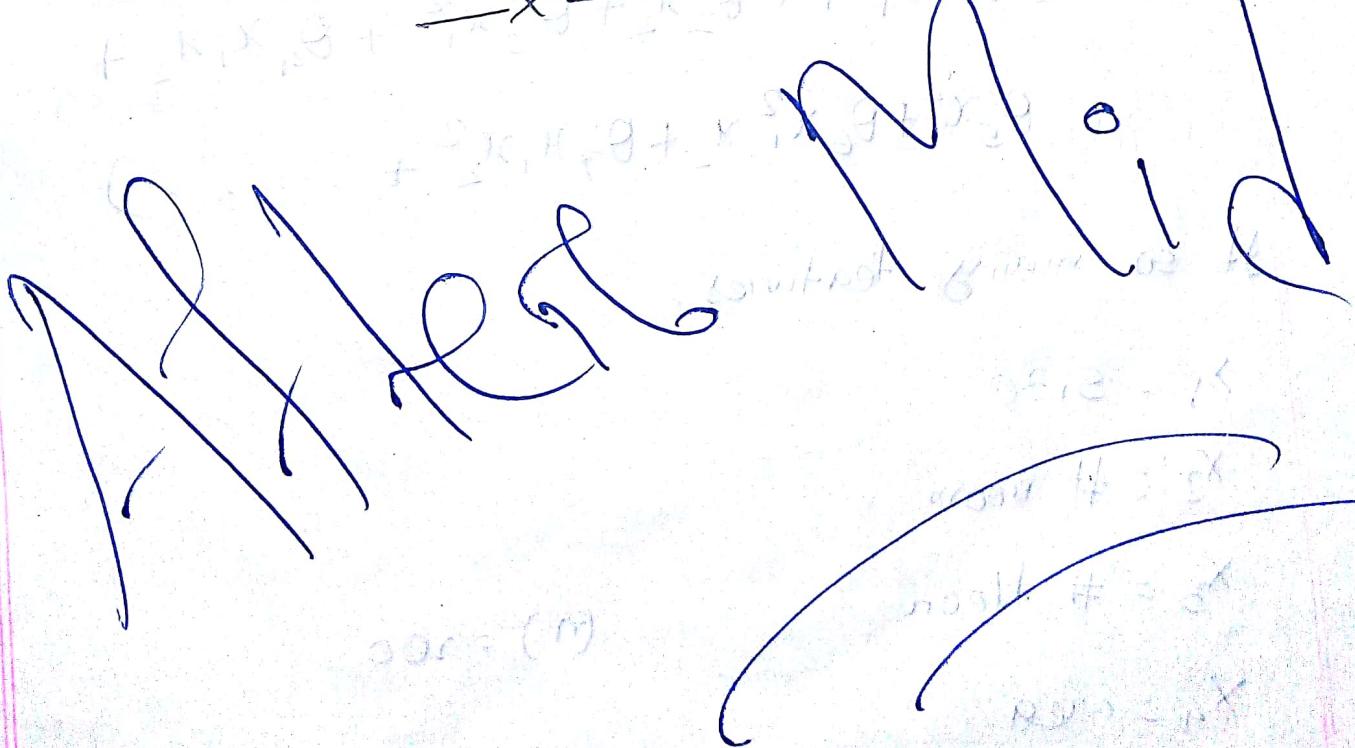
for i = 1 to n

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

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

Kaggle

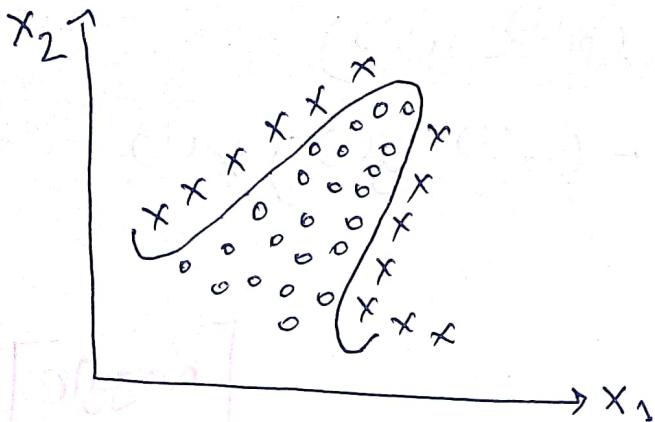
registration



24.06.18.

Neural Networks

* classification Problem



Q(2)

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

if so many features.

$$x_1 = \text{size}$$

$$x_2 = \# room$$

$$x_3 = \# floor$$

$$(m) = 100$$

$$x_4 = area$$

$$\dots$$

$$x_{100} = \dots$$

16

$$x_{100} = -$$

2^{nd} 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 $O(n^2)$

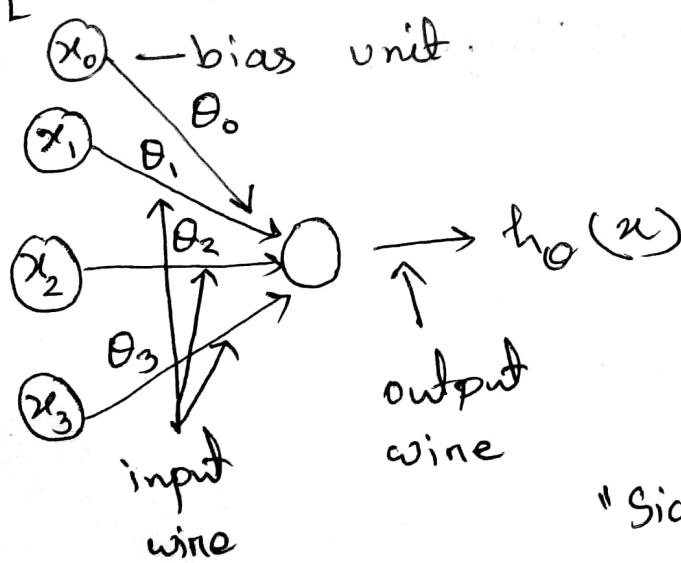
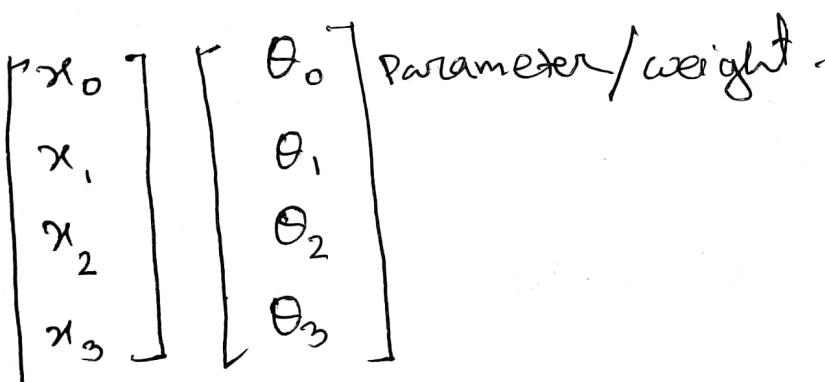
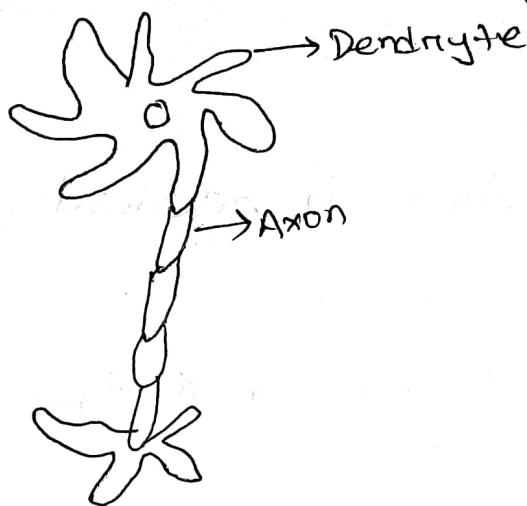
Gray image \rightarrow 8 byte

Color image \rightarrow 32 byte line

Artificial Neural Networks

26.06.18.

(for classification problem)

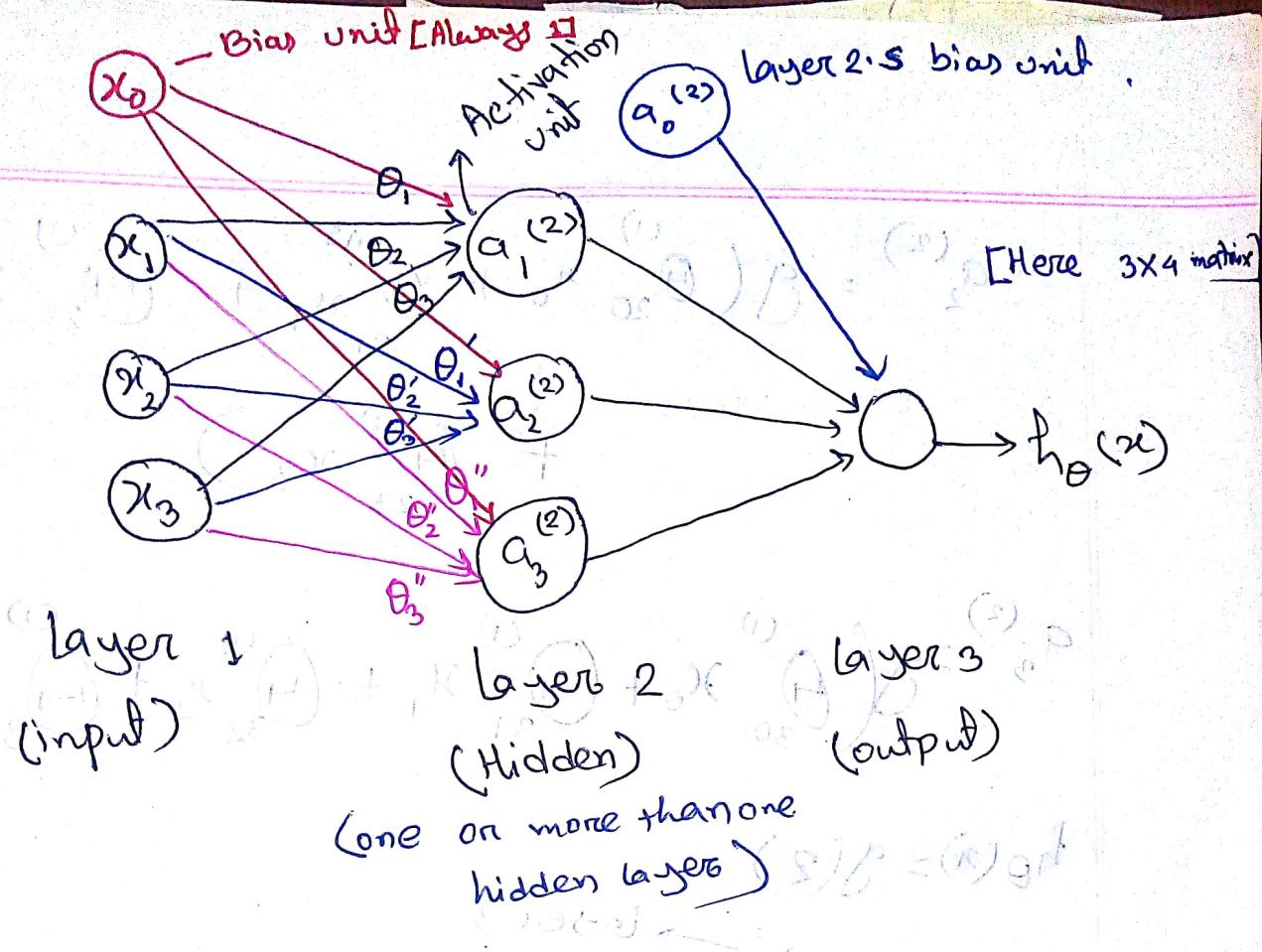


"Sigmoid Activation function"

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

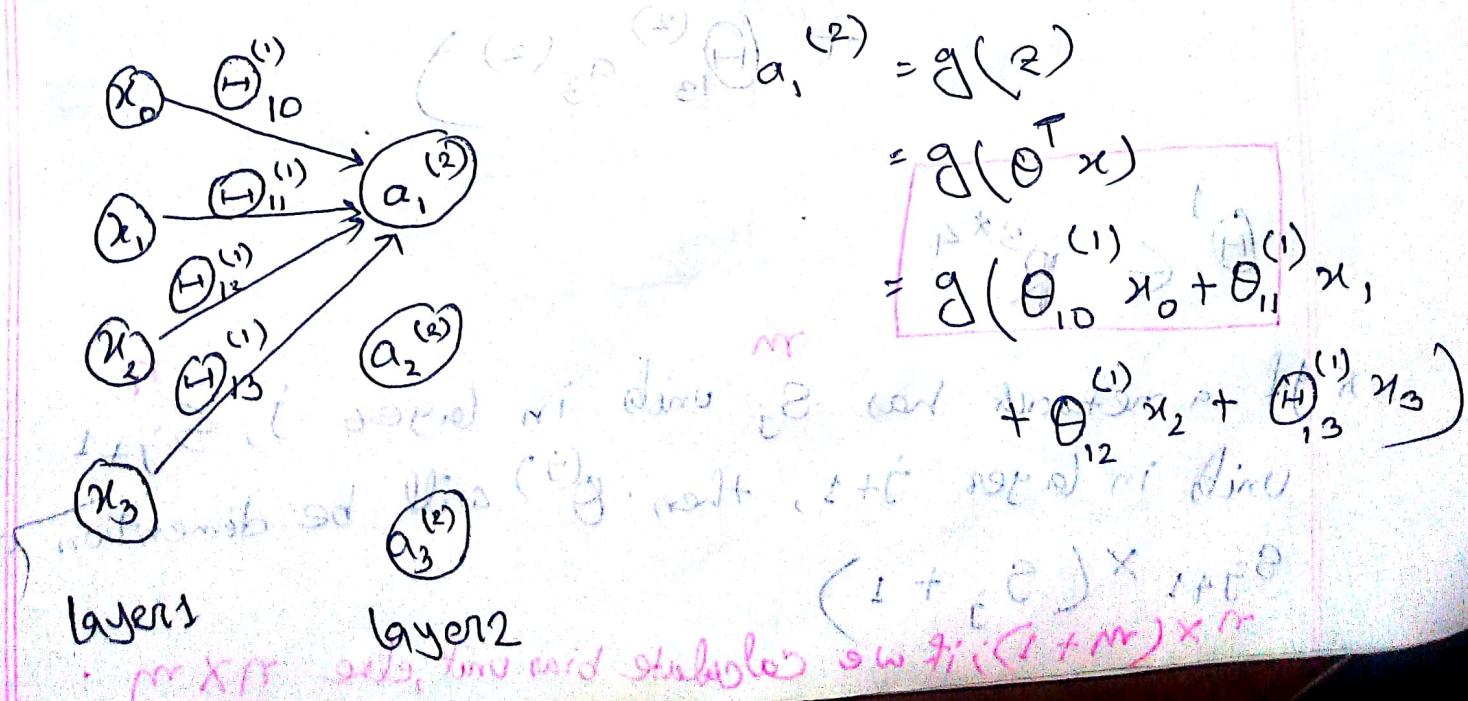
$$= \frac{1}{1 + e^{-\theta^T x}} \quad (\text{use linear order})$$

18



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

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



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

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

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

→ layer 3

$a_1^{(3)} = a_1^{(2)}$ → Unit 1

$$h_\theta(x) = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

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

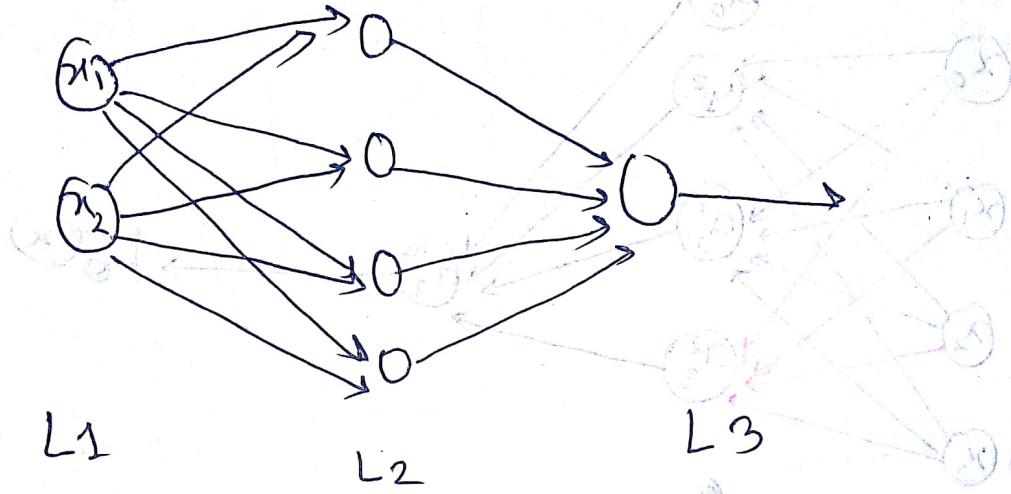
* If a network has s_j units in layer j , s_{j+1} units in layer $j+1$, then $\Theta^{(j)}$ will be dimension of $s_{j+1} \times (s_j + 1)$

$n \times (m+1)$; if we calculate bias unit, else $n \times m$.

20

6x4x3

i.e.



if we use bias unit, Dimension would be
 $(4 \times 3 + 1) \times 4 = 4 \times 13$

$$(4 \times (2+1))$$

$$(4 \times 3 + 1 \times 3 + 1 \times 3 + 1) \times 3 = 4 \times 3$$

$$(4 \times 3 + 1 \times 3 + 4 \times 2) \times 3 = 4 \times 13$$

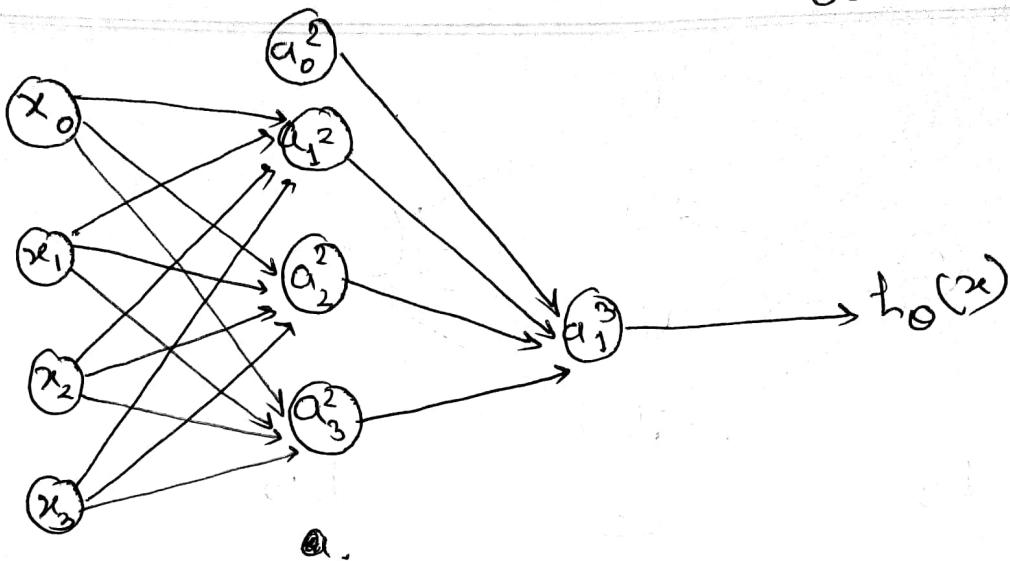
(Any)

$$(4 \times 3 + 1 \times 3 + 4 \times 2) \times 3 = 4 \times 13$$

(2) Bias frontmost output layer

$$(4 \times 3 + 1 \times 3 + 4 \times 2) \times 3 = 4 \times 13$$

$$(4 \times 3 + 1 \times 3 + 4 \times 2) \times 3 = 4 \times 13$$



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

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

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

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

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

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

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

3x4

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

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

For layer $j=2$ and k unit,

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

$$\frac{\underline{z}^{(3)}}{\mathbb{R}^3} = \textcircled{H}^{(2)} \underline{a}^{(2)}$$

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

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

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

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

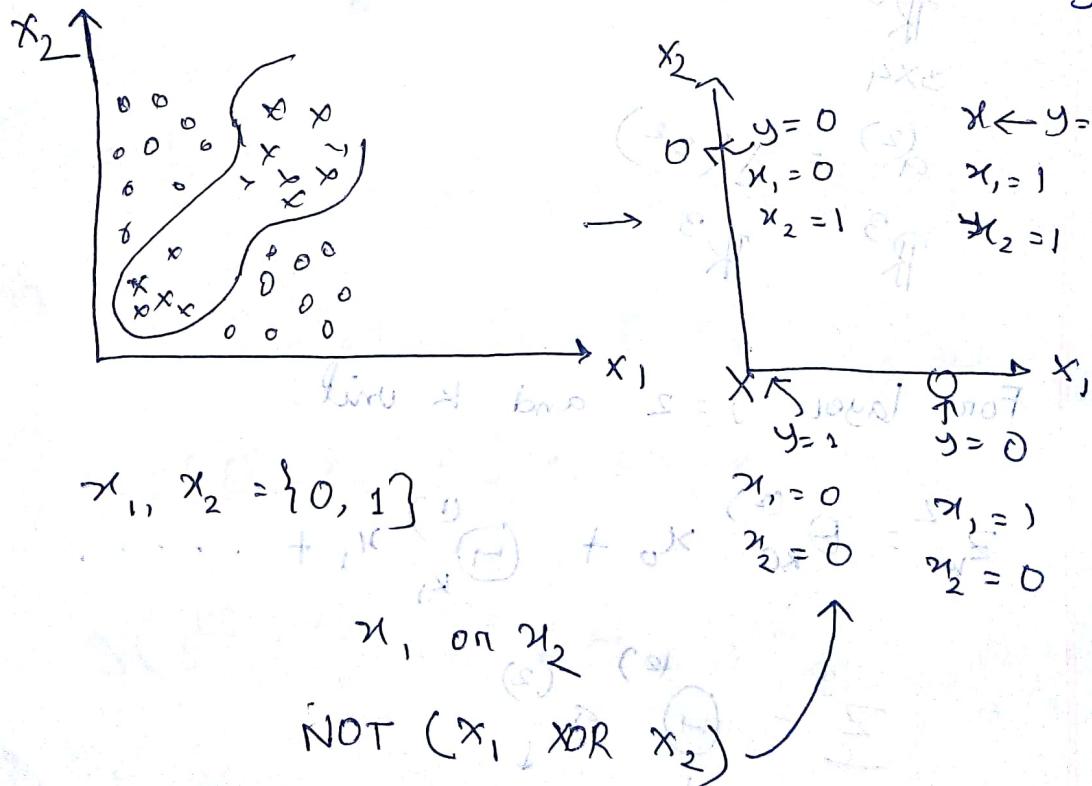
23

$$a_1^3 = g(z_1^3)$$

$$= \frac{1}{1+e^{-z_1^3}}$$

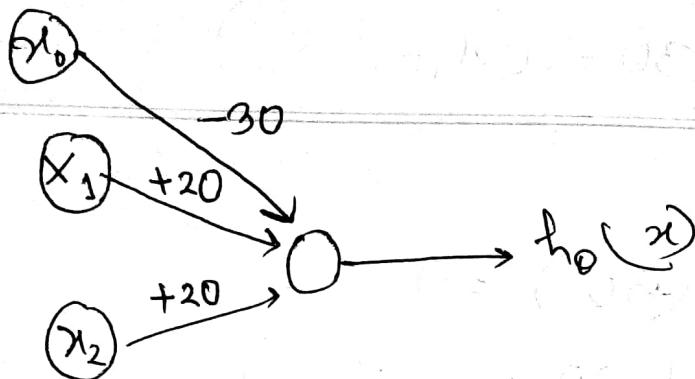
(logistic Regression)

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



AND, OR, NOT	x_1	x_2	$x_1 \oplus x_2$	$\neg(x_1 \oplus x_2)$
	0	0	0	1
	0	1	1	0
	1	0	1	0
	1	1	0	1

24



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

$$\begin{aligned} &= g(\theta_{00}^{(1)} x_0 + \theta_{10}^{(1)} x_1 + \theta_{20}^{(1)} x_2) \\ &= g(-30 + 20x_1 + 20x_2) \end{aligned}$$

$$\frac{x_1}{x_2}$$

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

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

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

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

$$g(2) = 1$$

x_1 AND x_2

$$z \geq 0$$

0

0

1

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

25

$$\frac{x_1}{0} \quad \frac{x_2}{0}$$

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

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

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

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

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

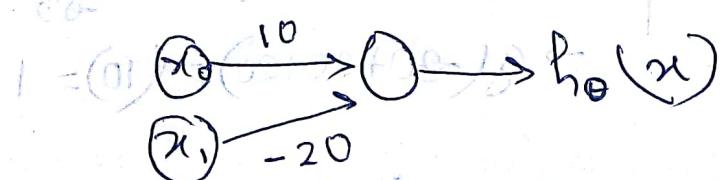
$$\frac{x_1}{0} \quad \frac{x_2}{0}$$

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

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

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

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



$$\frac{x_1}{0}$$

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

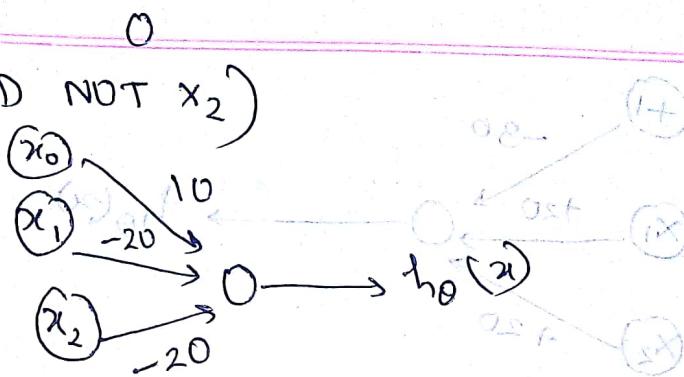
$$1 \rightarrow g(10 - 20) = 0$$

Assignment

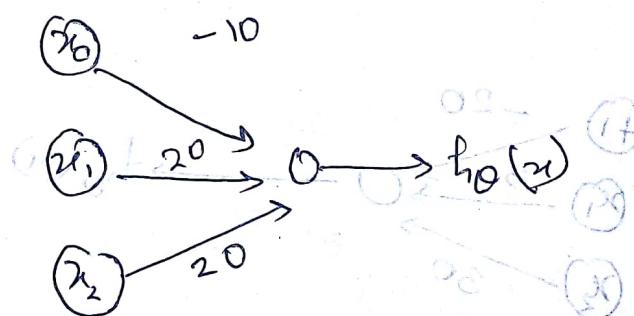
26

(NOT x_1 AND NOT x_2)

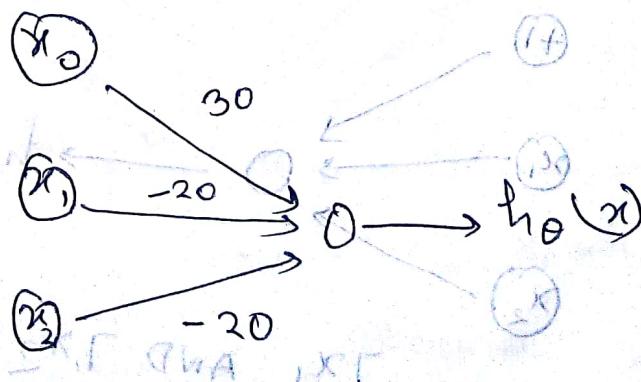
(a)



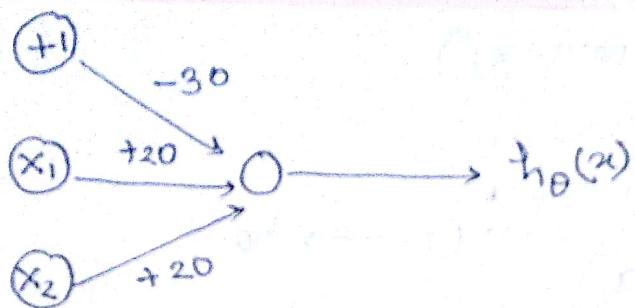
(b)



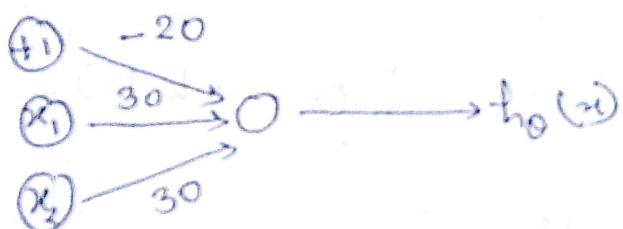
(c)



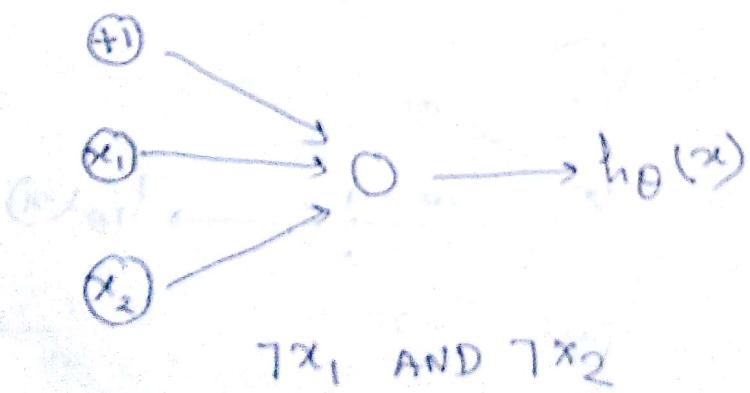
03.07.18. 27



AND

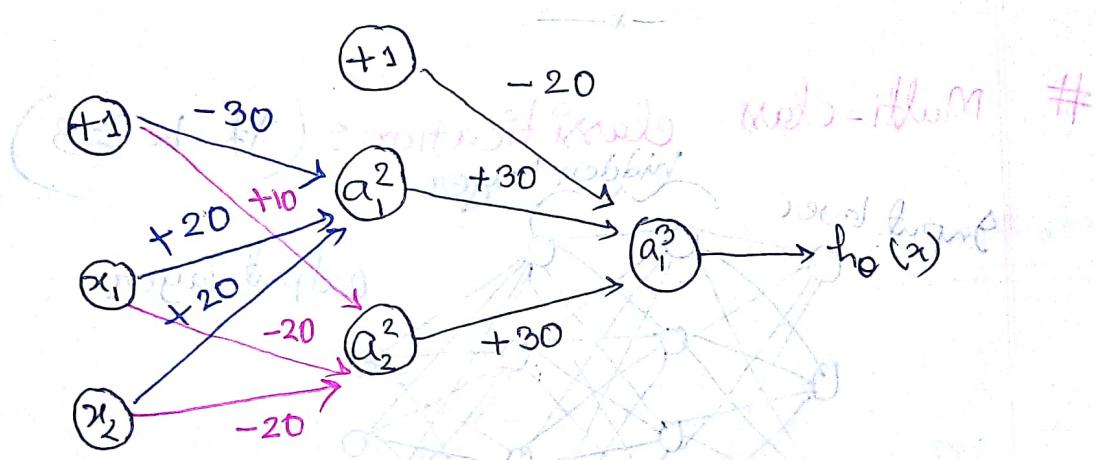
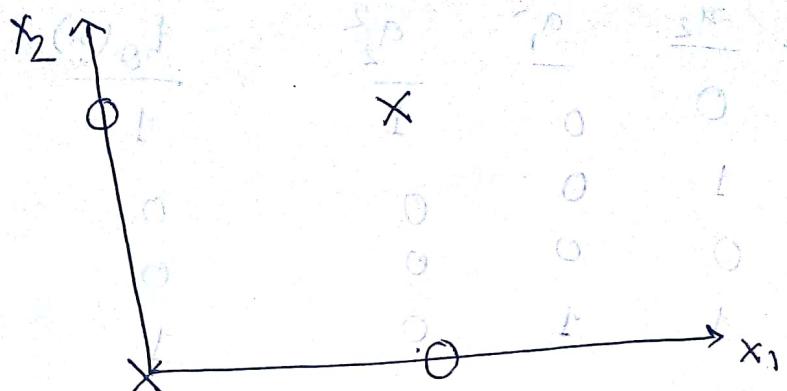


OR

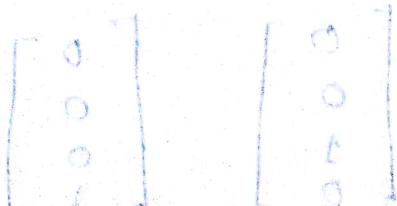
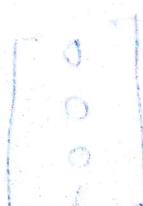


Assignment

28



(Forward propagation)



Input pattern



Result

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

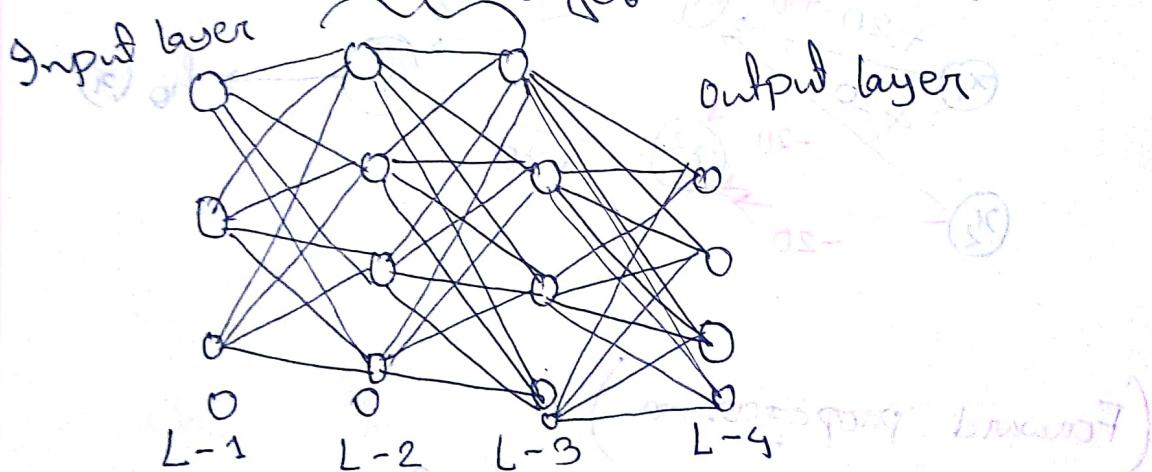
$$a_1^2 = 0 \text{ otherwise } 0$$

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

$$a_2^2 = 0 \text{ otherwise } 0$$

x_1	x_2	a_1^2	a_2^2	$\frac{h_0(x)}{1}$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

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



$$t = g_k = k + 1 \quad t = 3$$

$O_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$O_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$	$O_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$	$O_4 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$
car	truck	Human	motor cycle.

30

EX 10.80

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

training data consist of this batch of woh \leftarrow

$L = \# \text{ layers in network}$ \rightarrow pick first

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

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

$K = K$ class classification.

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

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

Binary classification

$$s_L = ? \quad 1$$

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

* If two class, then

$$(w_0 b + w_1 x_1 + \dots + w_n x_n) \geq 0 \quad \text{Binary classification.}$$

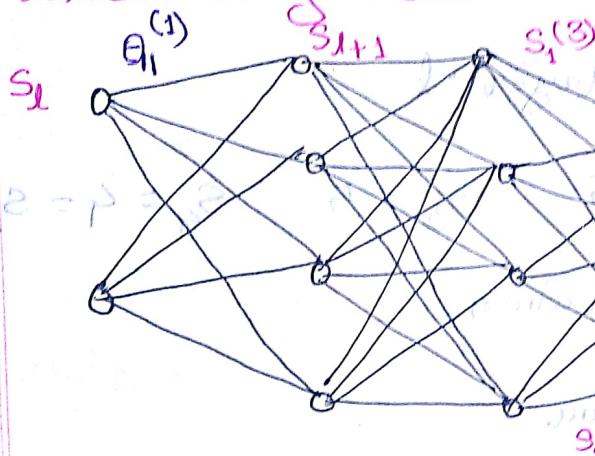
$$y \stackrel{\text{exit}}{\rightarrow} \frac{1}{1+e^{-\theta}} + [(w_0 + 1)b]$$

31

08.07.18

Today's Topic:

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



softmax L-1 (input) L-2

L-3

L-4

$k = 3$
multiclass
classification.

* Cost function:
 $J(\theta) = \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)})$

Logistic Regression:

soft. func. out θ

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

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

32

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

* Neural Network:

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

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 $\theta_{j|0}$ because

i starts with 1, so $\theta_{j|0}$ will be eliminated.

I assign j th layer to $j=0$: pointed
 j th unit. \rightarrow this logic does not

* Gradient Descent:

min

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

example

Forward propagation

$a^{(1)}$

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

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

$a^{(2)}$

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

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

$a^{(3)} = g(z^{(3)})$

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

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

* Gradient Computation: Back propagation algorithm:

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

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.

$$s^{(3)} \rightarrow s_1^{(3)} \quad s_2^{(3)} \quad s_3^{(3)} \quad s_4^{(3)}$$

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

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

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

, * = element wise
multiplication

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

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

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

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

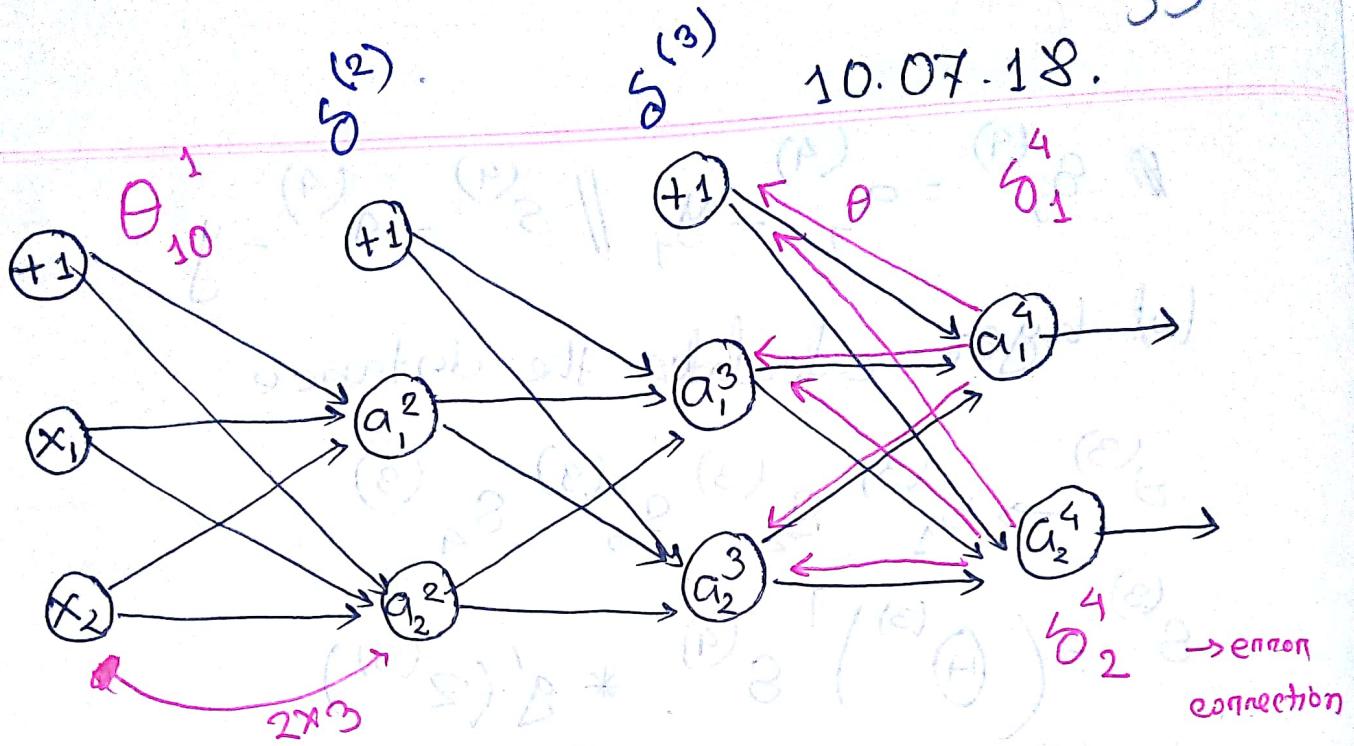
no error

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

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

35

10.07.18.

Cost Function: $J(\theta)$

$\min_{\theta} J(\theta) \quad | \quad \frac{\partial}{\partial \theta} J(\theta)$

layer 1 to layer 2
2's unit \times 1's unit + 1

$$(a_{j-1})^T * \dots * (a_1)^T * \delta_j^U = L \text{ layer, } j \text{ unit is error.}$$

Unit of
 $l+1$

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

BFSG

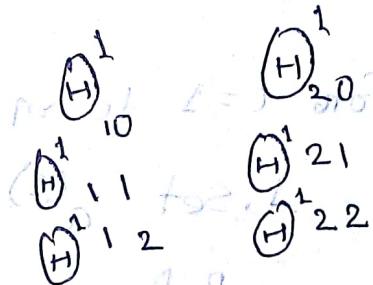
$$\frac{\partial}{\partial \Theta_{ij}} J(\Theta) \quad \text{takes}$$

↓ takes

$$J(\theta)$$

& partial derivative
o) input.

For layer -1



spotted ground antelope

$$s^{(4)} = a_j^{(4)} - y_j = a^{(4)} - y \quad \text{step 4}$$

$$g^{(3)} = \left(\begin{matrix} 0 \\ 0 \\ 1 \end{matrix} \right)^T g^{(4)} \neq g^{(1)} \left(\frac{1}{2} \right)$$

$$= (\theta^{(3)})^T g^{(4)} \cdot * a^{(3)} \cdot * (1 - a^{(3)}) \quad (\text{Back propagation})$$

$$\frac{\partial}{\partial \theta} j(\theta) = \hat{a}_j^{(l+1)} s_i^{(l)}$$

Back propagation algorithm:

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

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

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

- For $t = 1$ to m

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

- Perform forward propagation

to compute $a^{(l)}$ for

$l = 2, 3, \dots, L$

$$\begin{array}{ll} x_1 & a_1^{(1)} \\ x_2 & a_2^{(1)} \\ x_3 & a_3^{(1)} \end{array}$$

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

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

formula,

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

38

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

to update this below

Vector form

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

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

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

(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 is used as an "accumulator" to add up our values as we go along and eventually compute our partial derivative, this we get,

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

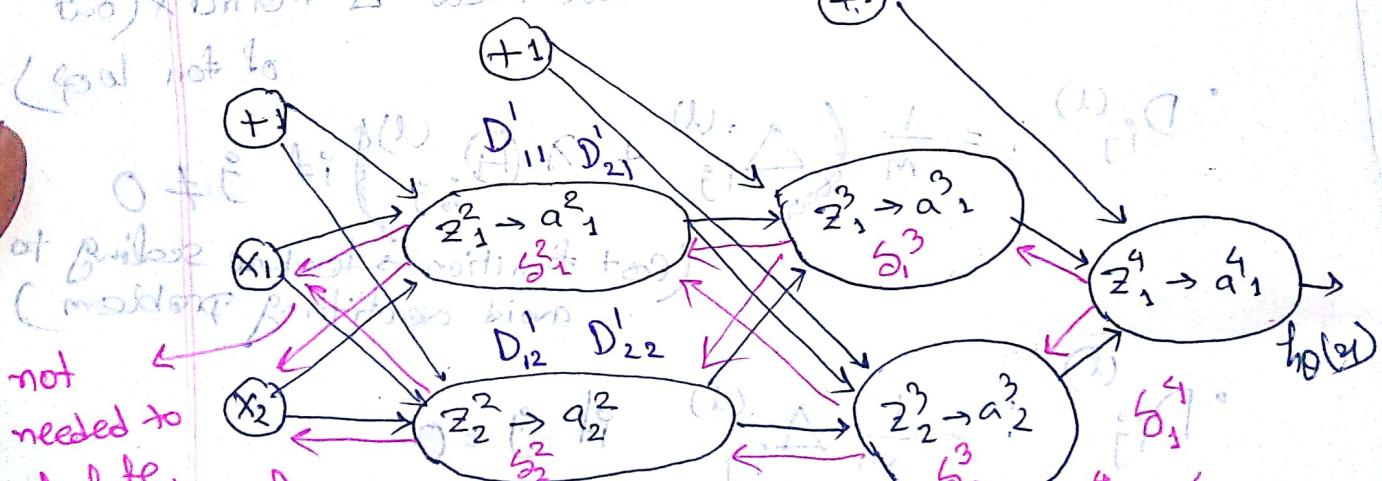
To reduce cost function, $J(\theta)$:

→ Gradient Descent

→ Conjugate GD

→ BFGS

→ L-BFGS



not needed to calculate, because original input has not errors.

$$z_1^3 = H_{10}^{(2)} \cdot x + H_{11}^{(2)} a_1^2 + H_{12}^{(2)} a_2^2$$

but at first time

get a number no. 91

Montage has grade of

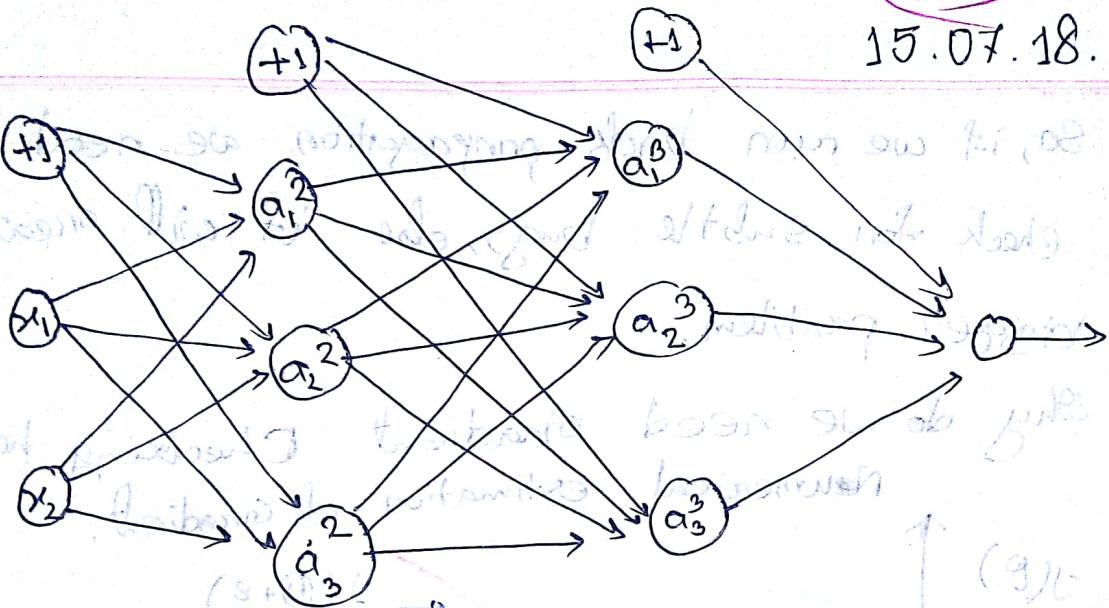
highest one steped

also see with derivative

$$\frac{\partial J}{\partial \theta} = (\theta) C \xrightarrow{S} \frac{\theta}{\theta_0}$$

(final)

15.07.18.



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

$$H^1$$

matrix
3x3

$$H^2$$

matrix
3x4

$$H^3$$

matrix
1x4



Gradient checking

We know how to do forward propagation
forward propagation
backward propagation

$$\frac{\partial J}{\partial \theta} = \frac{1}{n} \sum_{i=1}^n \frac{\partial J}{\partial \theta_i}$$

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.



$$\text{point } \theta_j = \theta_j \text{ given by } \frac{\partial}{\partial \theta_j} J(\theta) \quad \tan \theta = \frac{1}{x}$$

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

Parameter vector θ

allow sum of squares error metric to match response about all good

Unrolling Θ : $\Theta(\theta_0, \theta_1, \theta_2, \theta_3)$ to $\Theta(\theta)$

$\Theta^1 = 10 \times 11$

$\Theta^2 = 10 \times 11$

$\Theta^3 = 11 \times 11$

theta vectors = $\begin{bmatrix} \text{theta1}(:, :) & \text{theta2}(:, :) & \text{theta3}(:, :) \end{bmatrix}$

1-110 111-220 221-232

thetaVector[1]

Back to matrix

theta1 = reshape(thetaVector(1:110), 10, 11)

theta2 = reshape(thetaVector(111:220), 10, 11)

Parameter vector Θ

Let $\Theta \in \mathbb{R}^n$ (Θ is the unrolled version of $\Theta^1, \Theta^2, \dots, \Theta^n$)

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

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

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

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

$\{ \text{;} \} \text{ epsilon } \{ \text{;} \} \text{ constant}, \{ \text{;} \} 2\epsilon \text{ constant} = \text{constant}$

$$D_{1n} \frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_n - \epsilon)}{2\epsilon} \text{ (extension of 2nd)}$$

$\{ \text{;} \} \text{ constant } \{ \text{;} \} \text{ not constant} \text{ expand} = \text{start}$

for $i = 1 : m$

theta P = theta

theta Q(i) = theta P(i) + EPSILON

theta M = theta

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

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

end

check that,

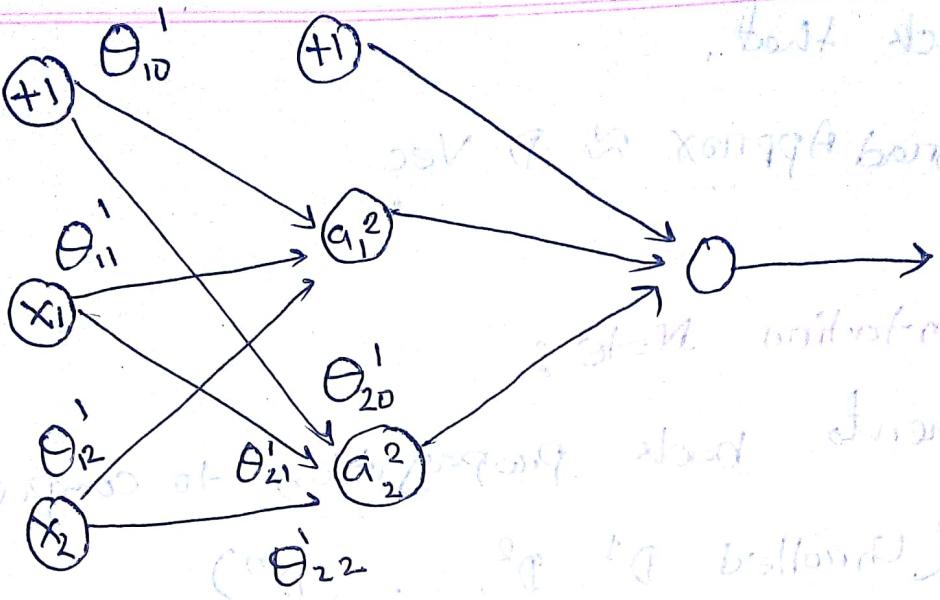
$$\text{grad Approx} \approx D \text{ Vec}$$

* 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 a 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.



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

$$a_2^2 = \theta_{20}^1 \theta_{21}^1 \theta_{22}^1 \text{ initialize to zero}$$

if next 0 present not do multiplying

$$a_1^2, a_2^2$$

$$\theta_{1x}^T \theta_{2x}^T$$

$$a_1^2 = a_2^2$$

and division ~~not~~ was done last time
X between \oplus two not

17.07.18.

Summary (all together):

First, pick a Network Architecture.

i.e. { 3 3 3 4

3 4 9 4
3 5 5 4

→ choose the layout of your Neural

Network.

↳ # hidden units in each layer

↳ # layers in total Network.

• # input units = dimension of features ($x^{(0)}$)

• # output units = # classes

• # hidden units / layer = ~~6*~~ ¹ ; Number

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

~~*2~~

Symmetry Breaking & Random Initialization.

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

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

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



Hidden Unit / layer = Usually more is better.

must balance with cost as it increase 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 $i = 1 \dots m$:

perform FP & BP using $(x^{(i)}, y^{(i)})$

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



Then compute $D^{(1)}$ 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.

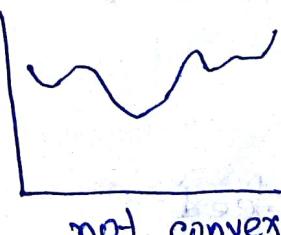
G. Then we use GD or a built-in optimization function to minimize the cost function with weight in theta.

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

It could be local minimum as well as global minimum.

a good local

optimal can get



not convex



convex

Convex

5960 //

~~CT-2~~

"Special class"

19.07.18.

Neural Networks:

Goal: not to be bio

Learn Protein Molecule

(G)C, sequence of proteins \rightarrow imagination

Human's Unsupervised learning \rightarrow Imagination
to predict the business plan of the company (new innovation)

9.8. own text contains at (G)C to this note

Architecture:

1. single \rightarrow layer feed forward neural network

2. multi "hidden" \rightarrow (D) deep part

How 3. Recurrent - all information is important

✓ Give \rightarrow How to choose hidden layer's units or nodes?

Data prepossessing \rightarrow fit data in model (G)C

Recurrent network:

* The Neuron:

1. Connection feed from last time

2. Adder

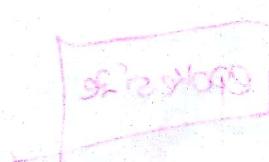
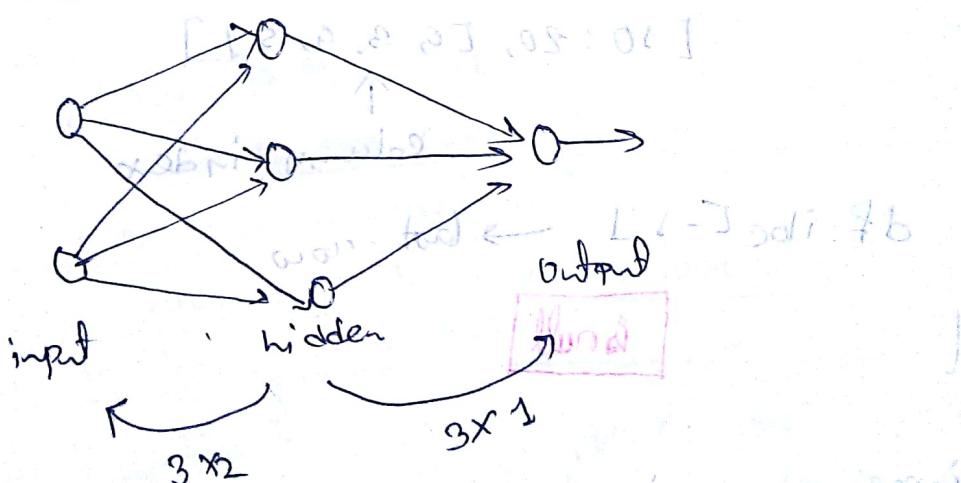
3. Activation

Dimensions of a NN:

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

* learning Algorithm: Backpropagation:

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



"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

high

isnull

linear regression.

epoch size

#

Association Rule Mining



Market Basket Analysis

1. Identify frequent item set → based on "support."

2. Based on frequent item set, discover association rule → satisfying "confidence."

i.e.: Consider the following transactions:

t₁: Bread, Butter, Milk

t₂: Bread, Butter, Coke, Diaper

t₃: Coke, Diaper, Butter

t₄: Bread, milk, Coke, Diaper

t₅: Coke, Butter, Milk

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

single item

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

(Support) → no support

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

$$\{milk\} = \frac{3}{5} = 0.6 \quad X$$

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

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

$\times \{Butter\}, \{coke\}$

Two items: $\{Butter, coke\}$

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

Association Rule: (confidence)

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

↑
if a person buys butter, then he also buys coke

(Apriori Algorithm) → (FP Growth Algorithm)
Best → frequent pattern (FP)

How many Butter, coke
Not combination (faster)

How many Butter

Min support & min confidence
depends on customer requirements.

$$18.0 = \frac{2}{5} = 0.4$$

could be

{ Butter → Diaper, coke }

{ coke → Butter, Diaper }

{ Diaper → Butter, coke }

* All possible combination depends on min. support.

Butter, Diaper, coke

Butter

→ X →

"Decision tree"

* Decision tree classification:

↓
Supervised learning

→ Phases

Input training set
(Training Data)

→ Inference →

to model creation

Test Data

Decision Tree algo.



model

(Decision tree model)

classify test data
(output)

i.e.: Example training data:

SL	color	Age	size	class
1	white	15	small	Tender
2	white	20	medium	Tender
3	gray	25	large	Mid range
4	gray	30	large	Mid range
5	Red	35	large	Mature
6	Red	36	large	Mature

1. calculate "Entropy"

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

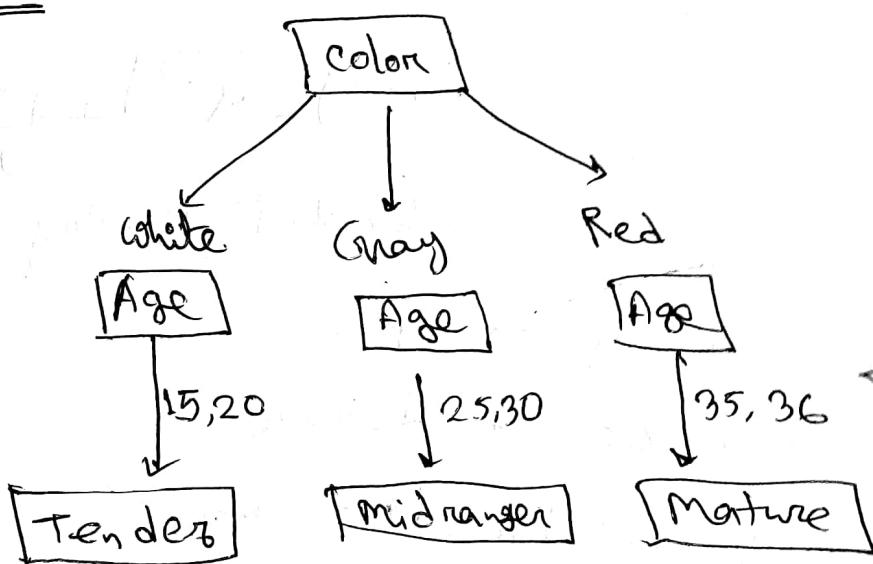
* Entropy:

→ inconsistency in data. How entropy is calculated
Higher the entropy, the higher the inconsistency.

Information Gain

→ To decide which attribute is best for splitting of the tree.

Trees



New Data :

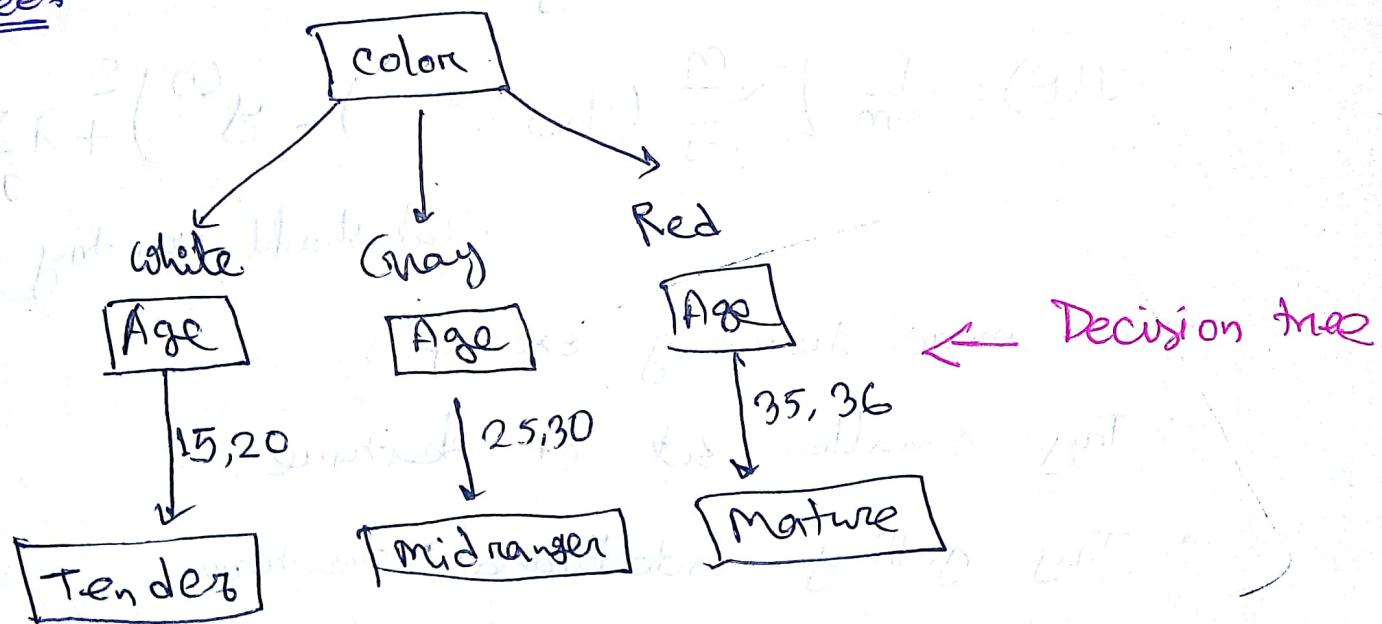
Red 32 (Not available)

— X —

Information Gain

→ To decide which attribute to choose as "root" of the tree based on information gain.

Trees



New Data

Red 32 (Not available).



22.07.18.

→ ML system Design

→ Understand multiple parts.

→ How to deal with skew data.

* Regularized Linear Regression (LR):

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (\hat{\theta}_0 + \theta_1 x^{(i)}) - y^{(i)} \right]^2 + \lambda \sum_{j=1}^m \theta_j^2$$

what should you try next?

→ Get more training examples.

→ Try smaller set of features.

→ Try getting additional features.

→ Try adding polynomial features

($x_1^2, x_2^2, x_1 x_2$, 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 of training example, we can split up the data into two sets,

- ① Training Set. (70%). Data may be used for training.
 - ② Test Set. (30%). Data used for testing.

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 Θ and minimize $J_{\text{train}}(\Theta)$ using training set.

2. Compute the test set errors $J_{\text{test}}(\Theta)$.

↳ LR

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

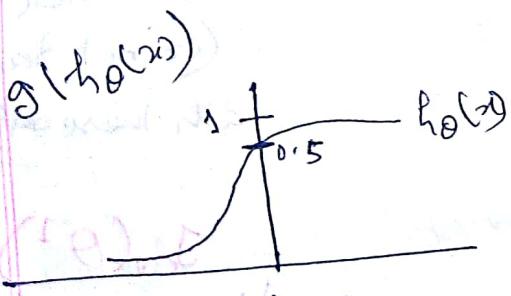
↳ classification Problem:

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

Missclassification errors,

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

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



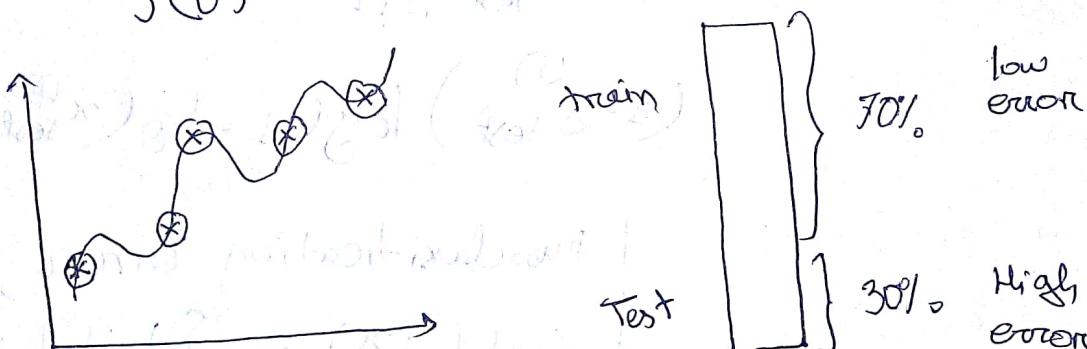
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

24.07.18.

$$\vec{\theta} \cdot \vec{x} = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

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

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



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% selected randomly

• cross validation set : 20% to validate

• Test set : 20%

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

1. Optimize the parameters in Θ 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 Θ_{test} .

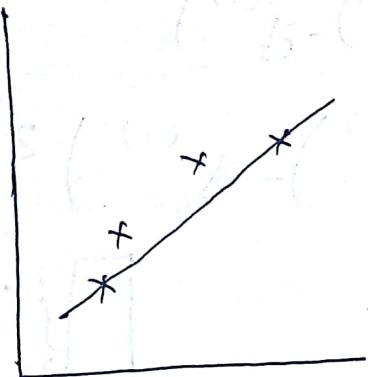
($d = \theta$ from polynomial with lower error)

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



29.07.18.

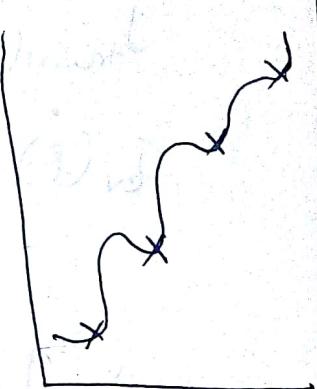
Diagnosing Bias Vs Variance:



$\theta_0 + \theta_1 x$
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
 $d = 4$

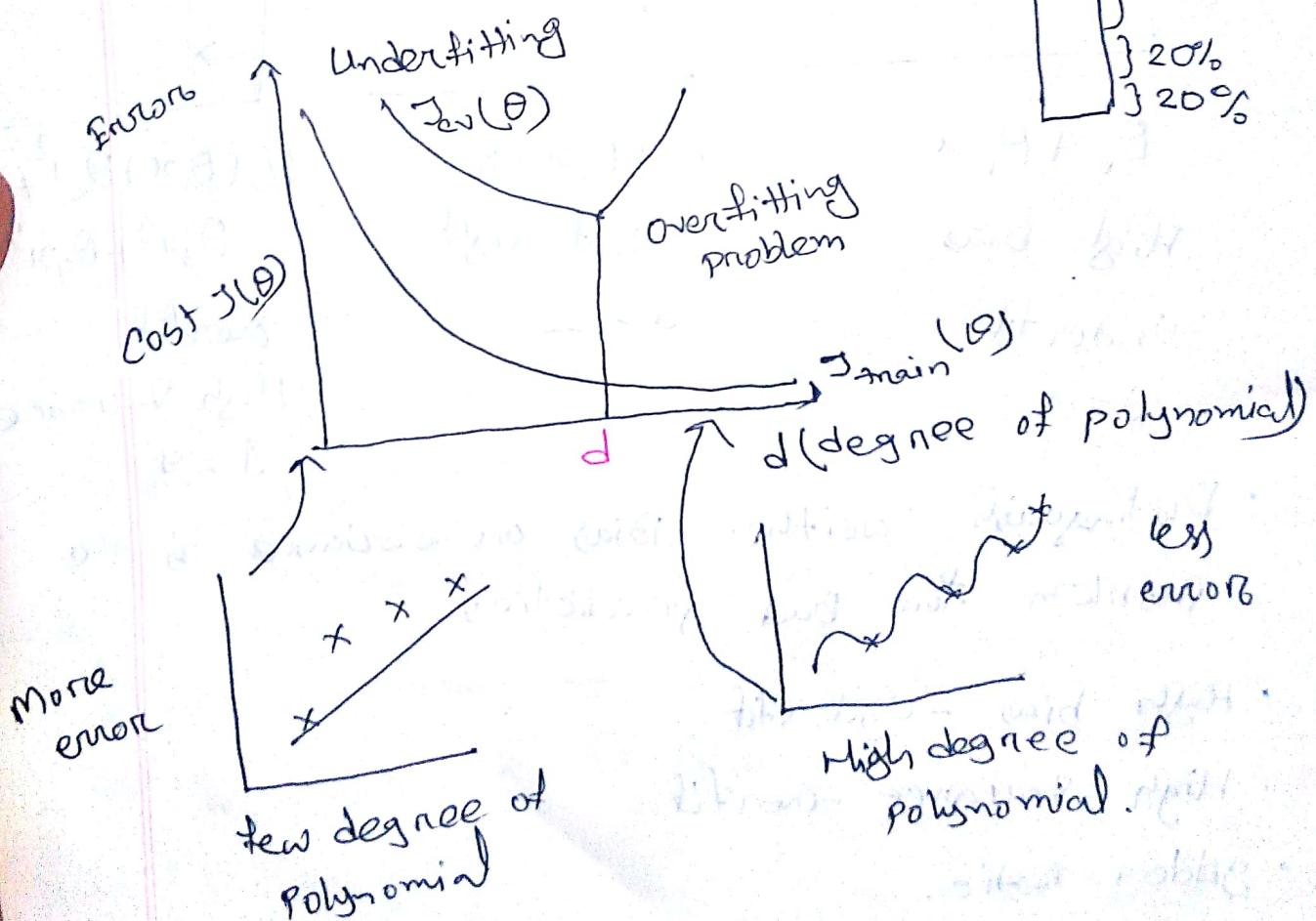
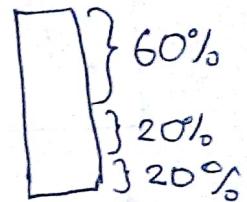
- Distinguish whether Bias or variance is the problem for bad prediction.
- High bias - Underfit
- High variance - Overfit
- golden ratio.

CV = Cross validation.

Bias/Variance:

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

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



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

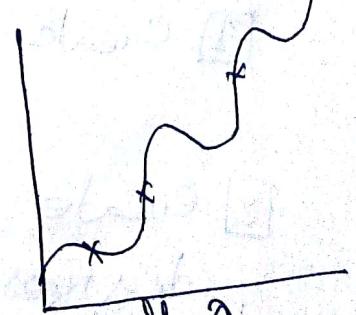
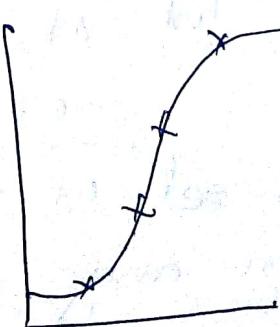
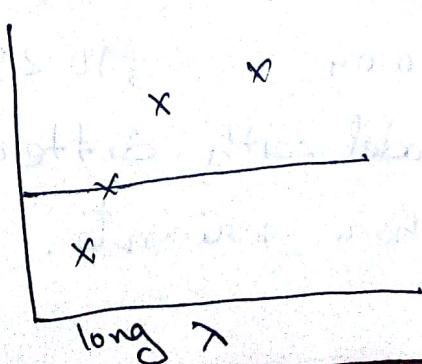
$J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$.

Regularization & Bias/Variance:

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

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

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



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

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

After selecting λ , How to Judge?

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

without regularization

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

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

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

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

λ , with less error will be selected.

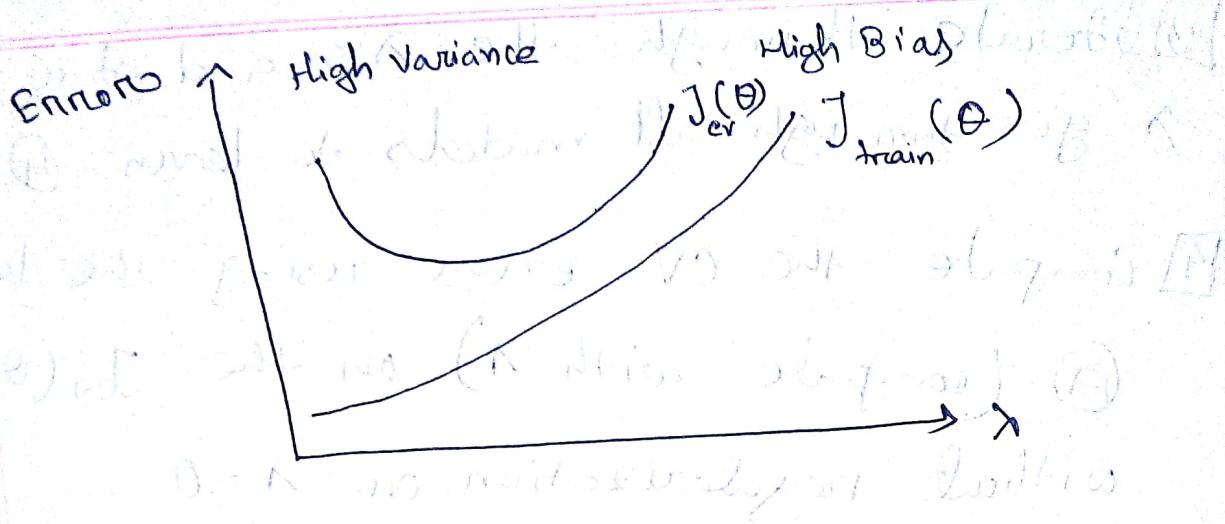
Choose the model and λ :

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

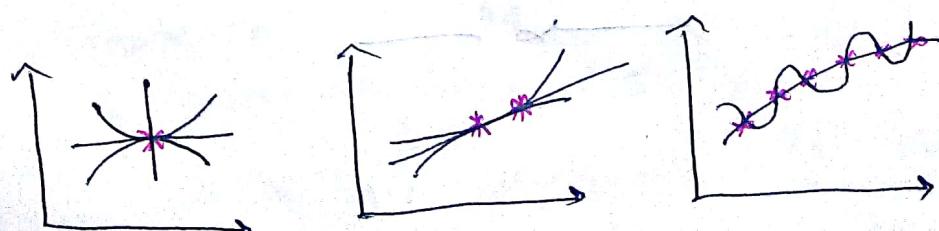
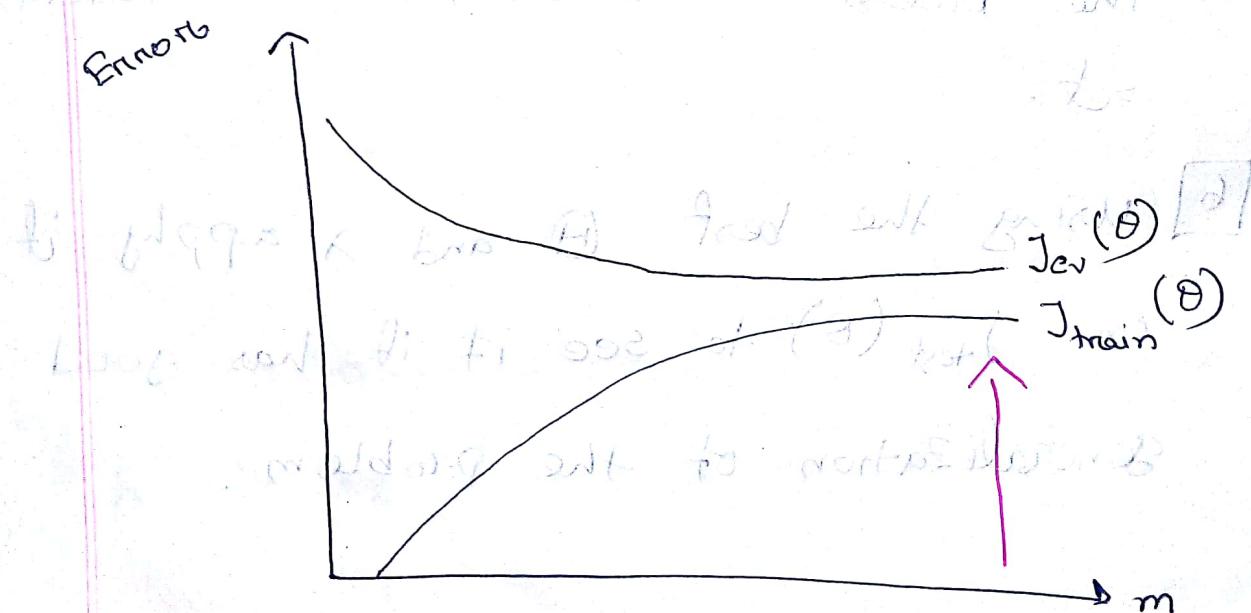
2. Create a set of model with different degrees or any other variants.

- 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 errors 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.

31.07.18



Learning Curves



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

→ Train algo. on a very few data, it will

↳ easily have 0 error.

↳ can find a quadratic function that touches all point.

→ Training set larger,

↳ Error increases

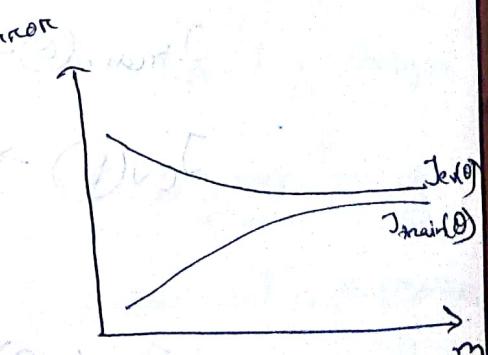
→ Errors value ^{will} plateau out after a certain m.

Experiencing High Bias:

→ Lower training set size,

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

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



→ Higher training set size,

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

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

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

* If learning alg.: 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,

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

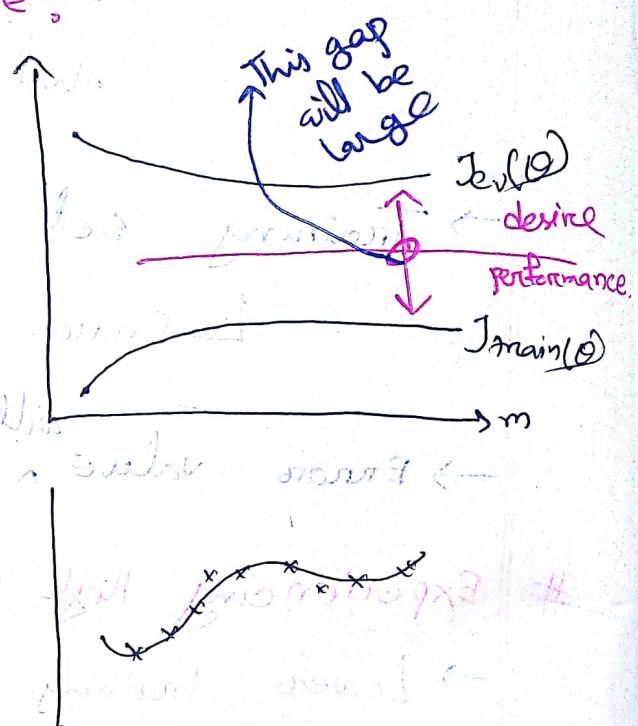
$$J_{\text{cv}}(\theta) \rightarrow \text{Decrease}$$

(But not low)

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

$J_{\text{cv}}(\theta) \rightarrow \text{Decrease}$ 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 $\hat{y}_0(x) = \begin{cases} 1 & \text{if cancer} \\ 0 & \text{if not} \end{cases}$

Your algorithm got 99% accuracy

0% Error.

0.05% Error

function $\text{pred}(x)$

$y=0$ # ignore x

last loop and result shows 99.5% Accuracy

so your model exhibit 0.05% Error.

using this we can do next and so on

Error Matrix:

Precision / Recall

$Y=1$. cancer

$$h_0(x) = 1$$

Actual class

		1	0
Predicted class	1	True Positive	False Positive
	0	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.

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

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

Algo. \rightarrow high P \rightarrow R

* If 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

$$0 \leq h_0(x) \leq 1$$

Predict 1, if $h_0(x) \geq 0$

Predict 0, if $h_0(x) < 0.5$

Predict $Y=1$, only when very confident

Predict, $Y=1$ if $h_0(x) \geq 0.7$.

" $Y=0$ if $h_0(x) < 0.7$

(Higher P, Lower R)

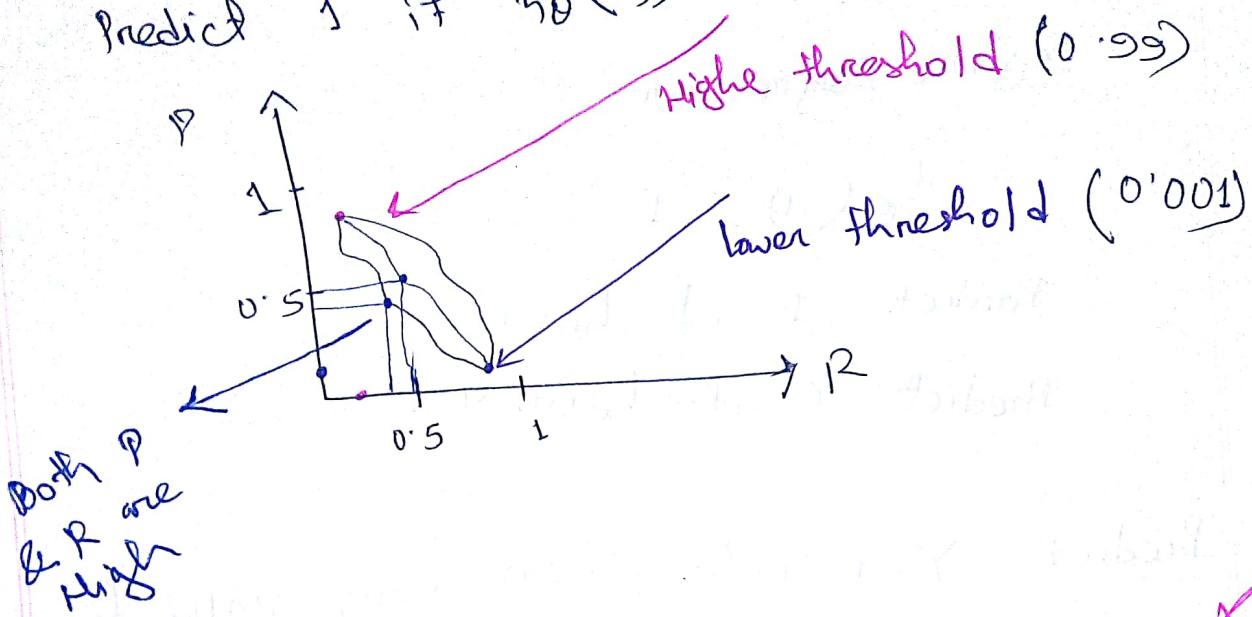
* Want to avoid missing.

Predict 1, if $h_0(x) \geq 0.3$

" 0, if $h_0(x) < 0.3$

(Lower P, Higher R)

Predict \rightarrow if $h_0(x) >$ threshold.



F₁ Score (F score) %

$$\frac{P+R}{2} \quad \frac{F_1}{0.44}$$

	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	0.505	0.0392

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

$R=0$ ~~P=1~~ $P \approx 1$

If $Y=0$

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

* what will be the approach in case of choosing algo of worse 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]

— X —

04.09.18

Support Vectors Machine:

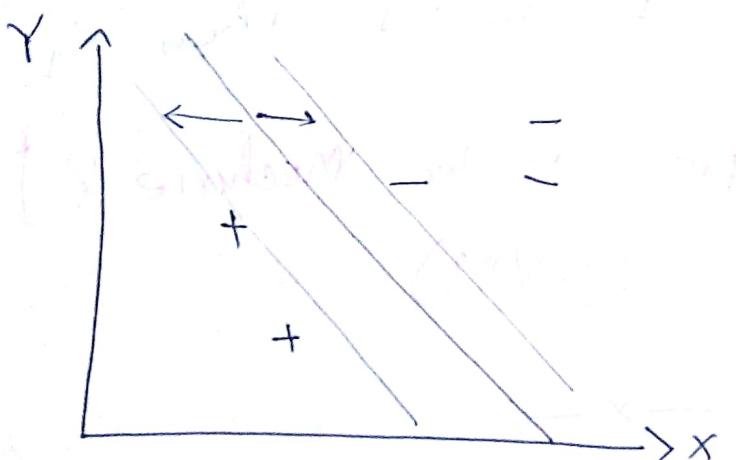
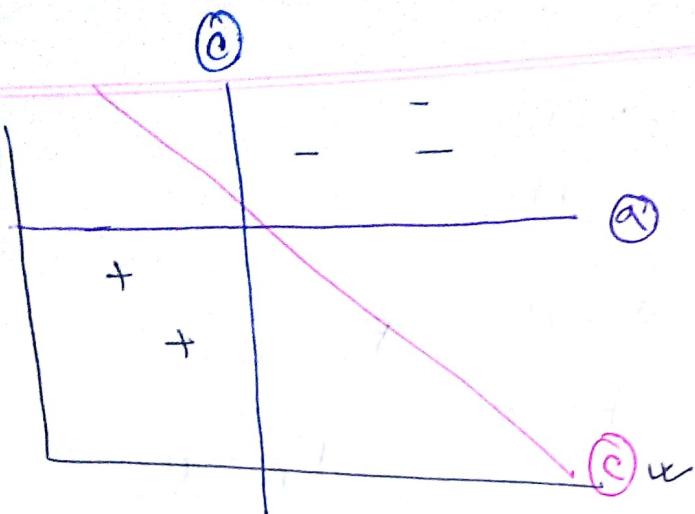
"It's not who has the best algo. that wins."

It's who has the most data."

When TRUE, when FALSE?

i.e.: For breakfast I ate — eggs.

to, two, two ?



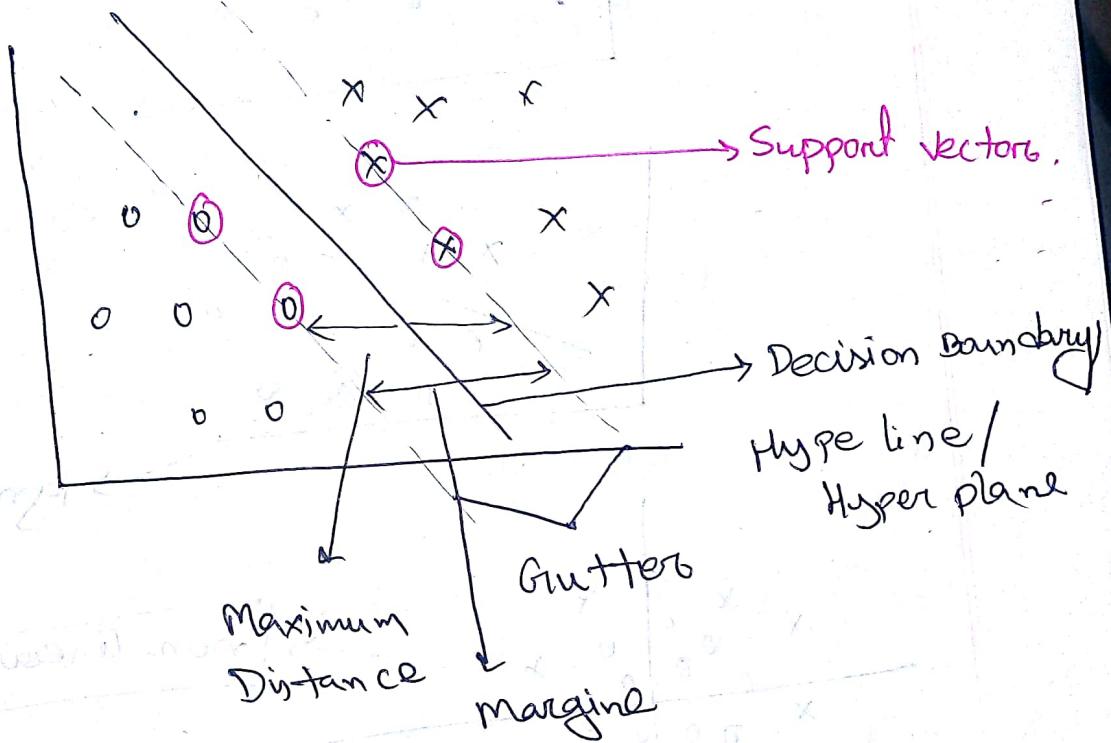
SVM \rightarrow ML. Algo. (Supervised)

\searrow Numerical classifiers
regression/classifiers.

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

SVM

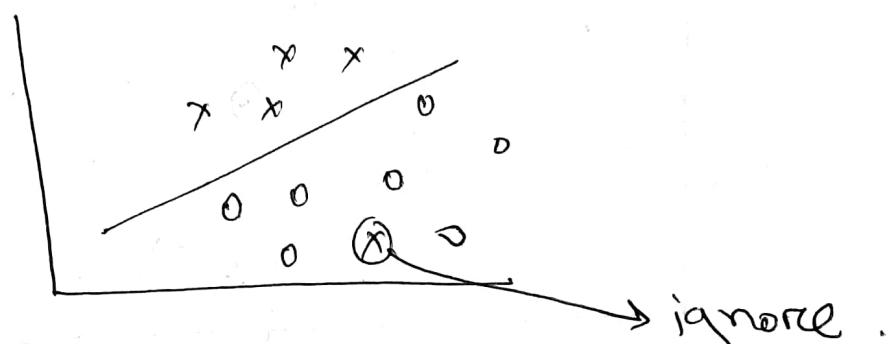
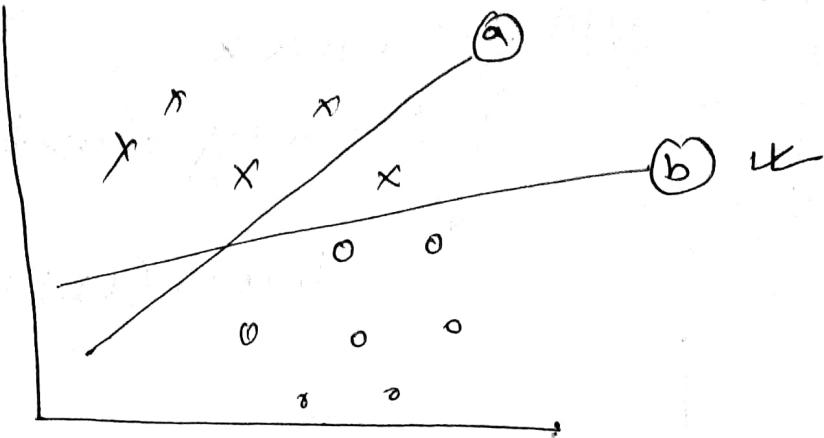
- convex
- Binary classifiers
- non-linear (KERNEL)



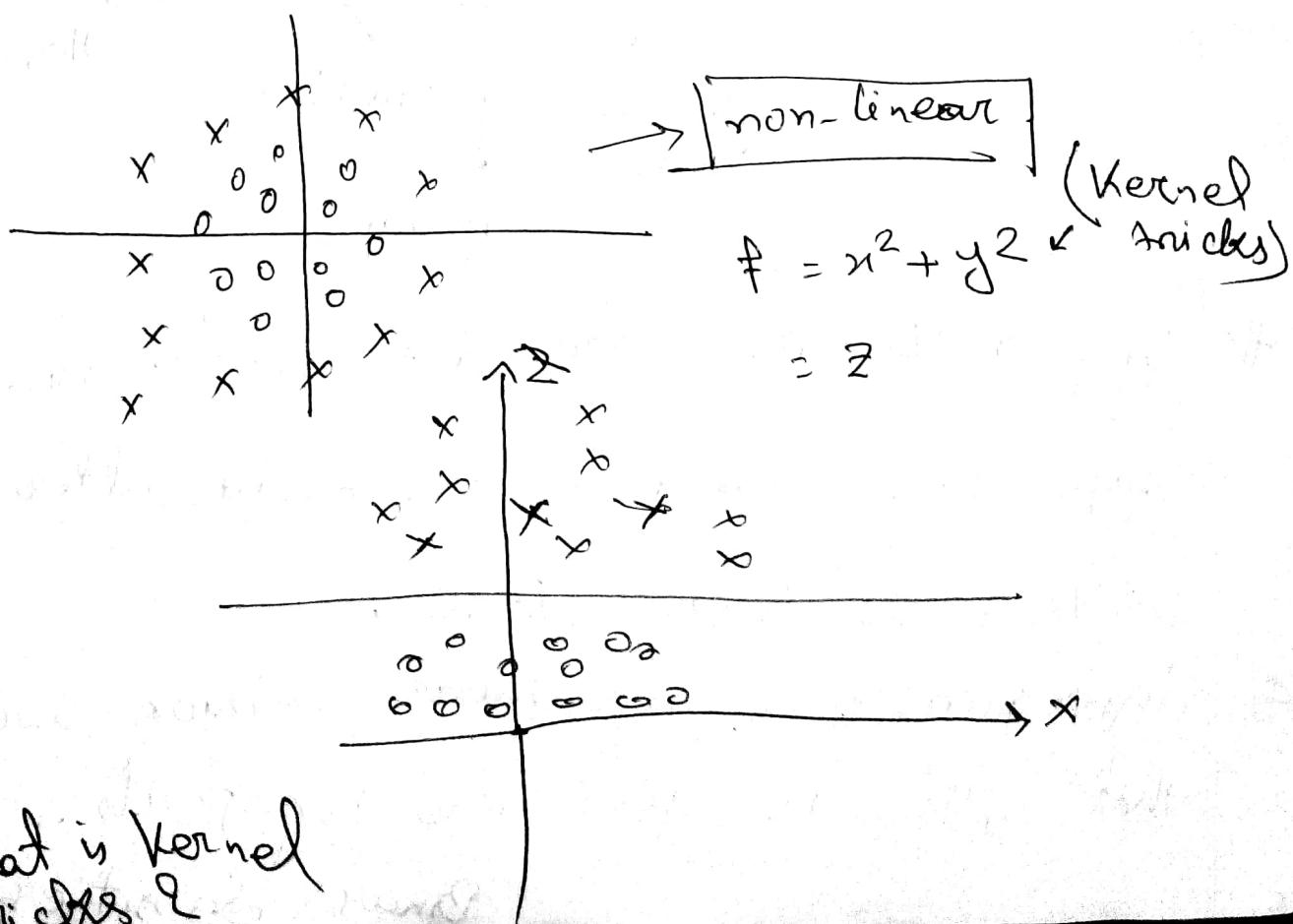
Support Vectors: are data points nearest to hyper plane. If removed, would alter the position of hyper plane.

Hyper plane: is a linear decision surface that splits the space into two parts.

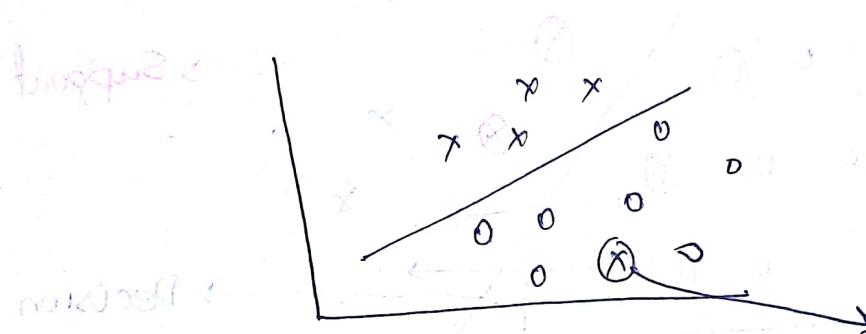
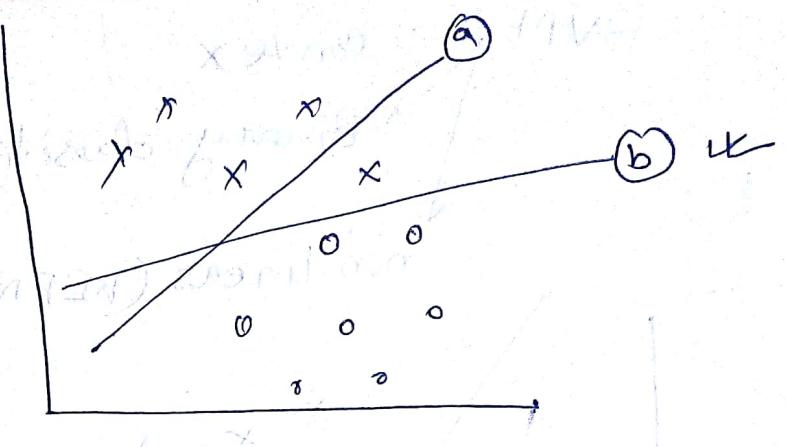
Binary classifiers



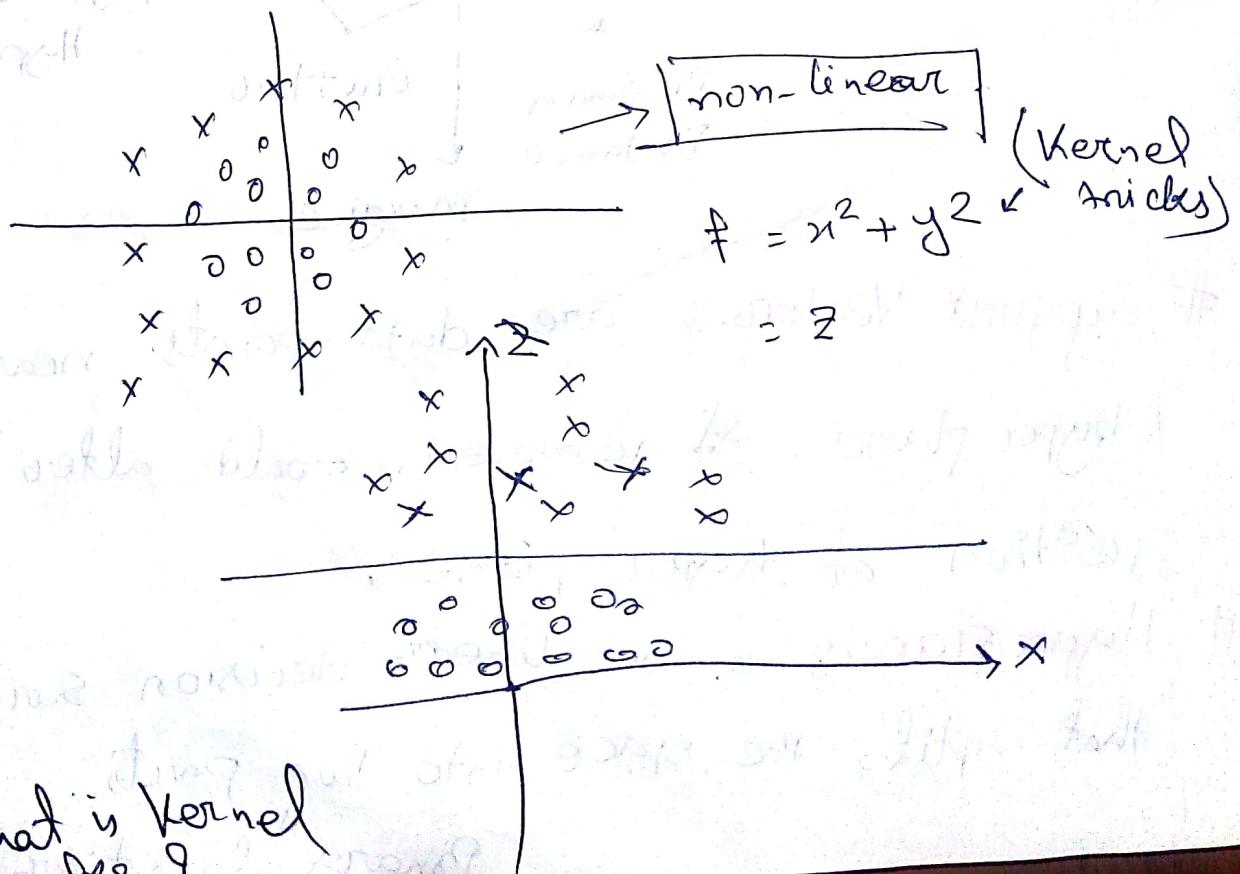
ignore .



* what is Kernel tricks?



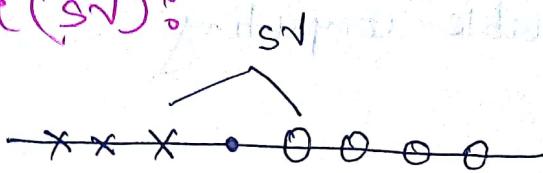
ignore.



* What is Kernel
tricks?

* # of Support Vectors (SV) :

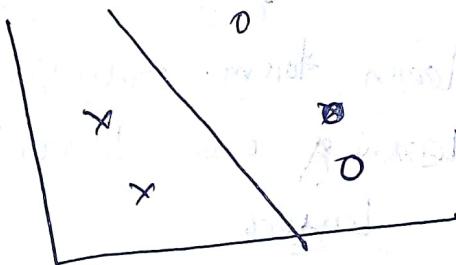
1 Dimension :



$$\# \text{ of SV} = 2$$

D : point.

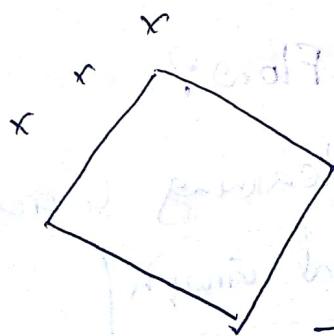
2 Dimension :



$$\# \text{ of SV} = 2/3/4$$

line.

3 Dimension :



* How to select SV? [Convex rule] \rightarrow H.G. (AT)
Jessica Noelle [MIT]
machines
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's

(i) construction phrase

(ii) Execution phrase

Everything is an object in TensorFlow.

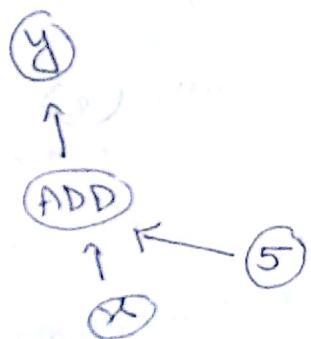
* What is a Tensor?

Tensor data structure to represent all data.

R

* Basic of TensorFlow:

(1) Variable



* linear Algebra:

(1) Dimensionality

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

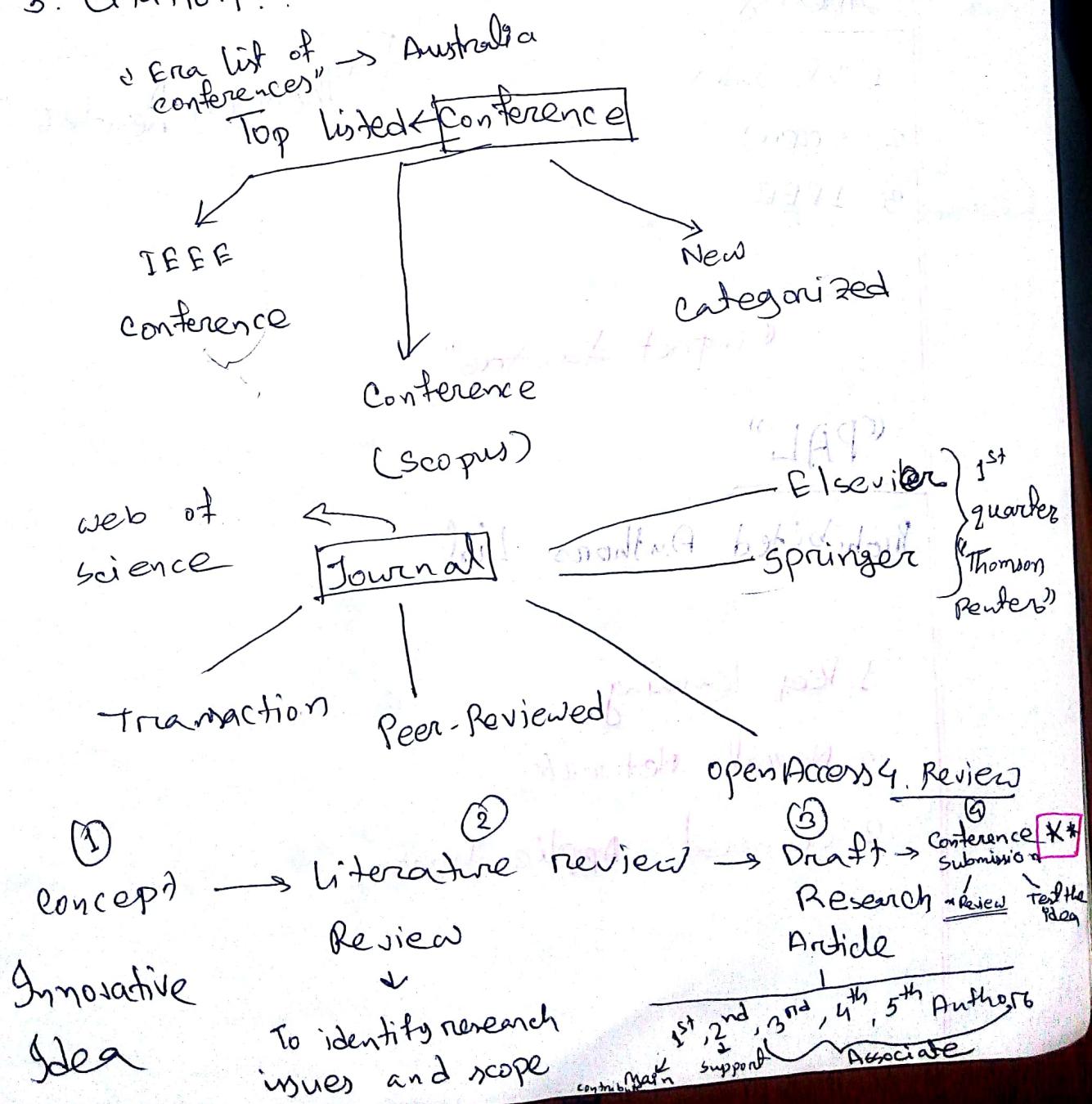
(1)

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

linear Regression:

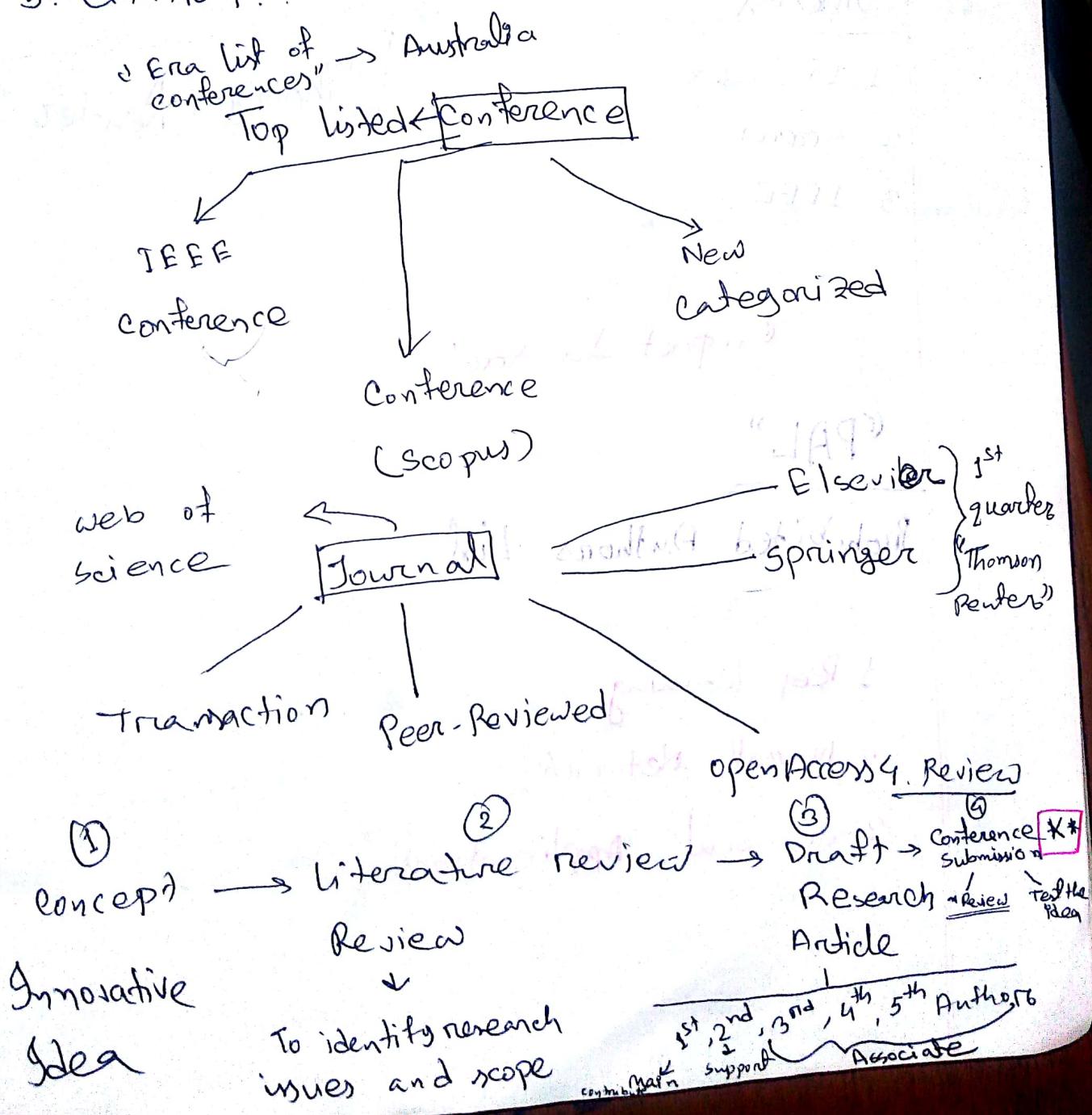
Research Discussion:

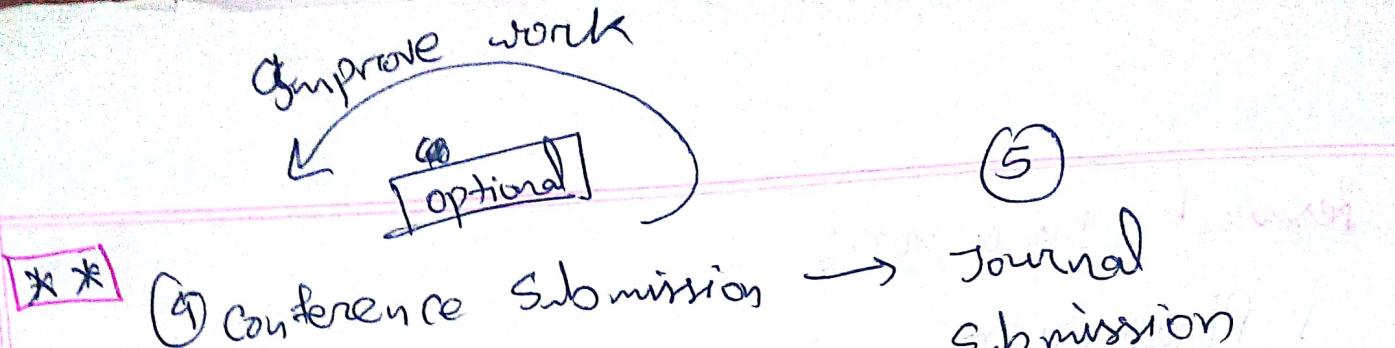
1. Plagiarism
2. Paraphrasing
3. Citation .



Research Discussion:

1. Plagiarism
2. Paraphrasing
3. Citation.





Indexing

1. ISI Index
2. Scopus
3. IEEE

"Thomson Reuters"

"impact factor"

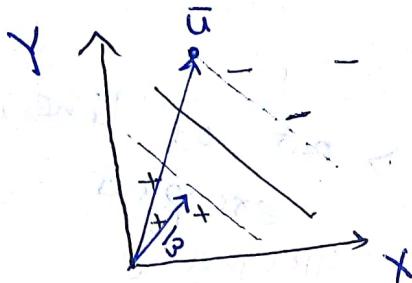
"PAL"

Prohibited Authors List

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

09.09.18.

Support Vector Machine:



vector is perpendicular
with my per plane,

\bar{w} = unknown point / new
example.

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

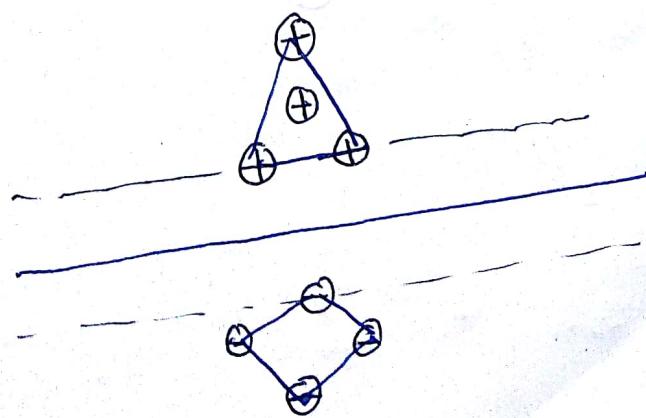
then +

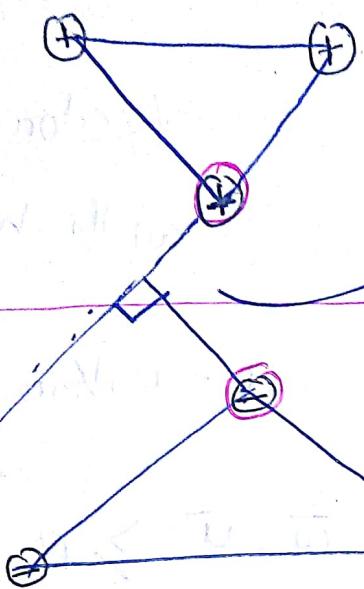
[b, c ,
constant]

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

α = Supportiveness

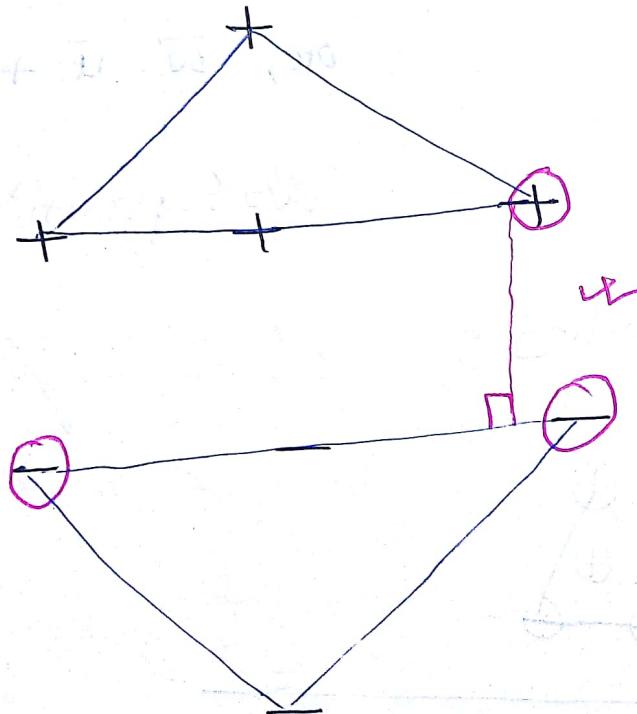
Convex hull Method:





 ||
 =
 ✓
 as (iii)
 +
 -
 2 points / SV

out of line, so not
 excepted
 hyper plane

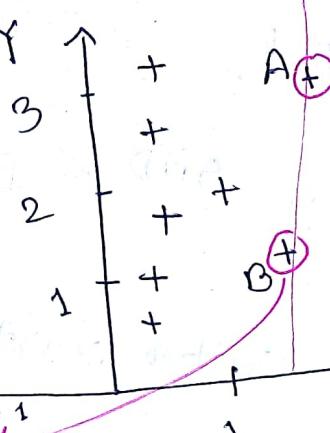


3 points / SV

calculating b , \bar{w}

$$\|\bar{w}\|$$

$$= 2$$



$w(\bar{w}) = [-2 \ 0]$
1. Draw the decision boundary. (DB)

2. Write the equation for the DB.

$$x=2$$

[2 dimensional vector]
as here two features

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

$$\bar{w} \cdot \vec{x} + b = 0 \quad \begin{bmatrix} x \\ y \end{bmatrix}$$

$$x_1 \theta_1 + x_2 \theta_2 + \theta_0 = 0$$

$$\theta^T x$$

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

$$\begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + (-\theta_0) = 0$$

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

4. Scale equation so that,

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

equivalently, use the gutter constraint.

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

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 $S \cup D$

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

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

$$2 - 2.5 = -1$$

$$-0.5 = -1$$

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

$$\Rightarrow 0.5 C = -1 \quad \left| \begin{array}{l} \vec{\omega} = \vec{\omega} \cdot C \\ = [1 \ -2 \ 0] \\ b = +3 \end{array} \right.$$

$$\Rightarrow C = -2$$

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

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

→ if margin decreases,
w will increase

else, decrease

by changing dimension or by increasing margin
or decreasing weight (which is increased or decreased).

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

$\|w\|$ is another term for norm of w

$$\text{margin width, } = \frac{2}{\sqrt{w^T w}}$$

or $\sqrt{w^T w}$ is called norm of w

$$= \sqrt{w^T w}$$

H.W →

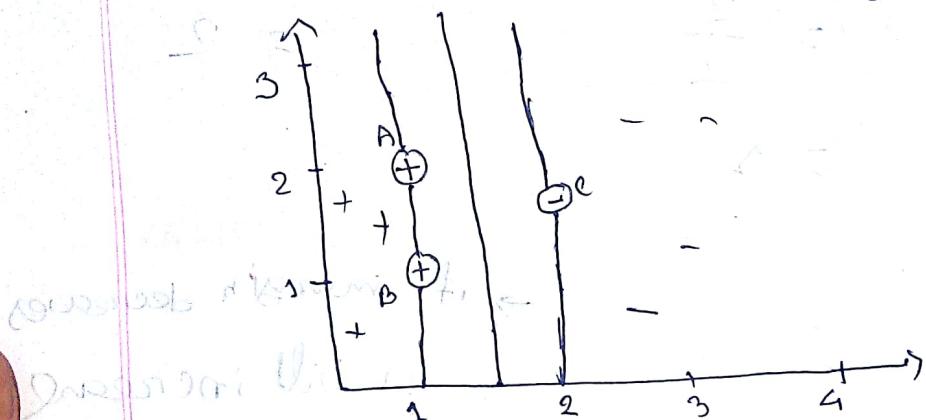
Winston
MIT professor

SVM

11.09.18

Calculating α (Lagrange multipliers) /

~~Supportiveness value~~ \downarrow
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 vectors, $\alpha = 0$

Otherwise, the point is a SV and

lies on the gutter and $\alpha > 0$

[Positive Supportiveness]

$$2. \sum_{i=1}^m \alpha_i y_i = 0$$

Support
vectors

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

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

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

$$\Rightarrow \bar{w} = \sum_P \omega_P \bar{x}_P - \sum_N \omega_N \bar{x}_N$$

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

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

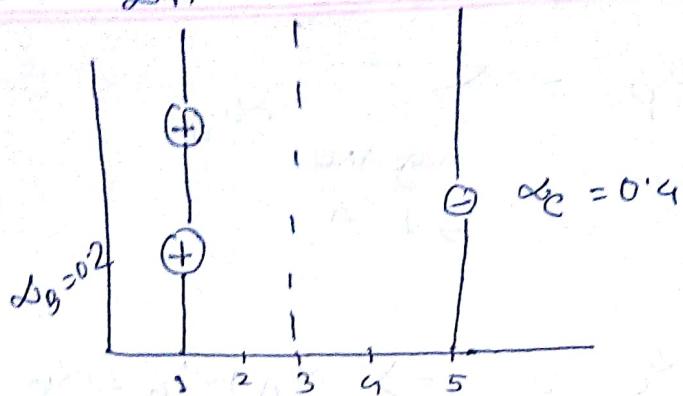
$$\Rightarrow -2 = \omega_A \cdot 3 + \omega_B \cdot 1 - \omega_C \cdot 2 \quad \text{--- (3)}$$

$$\omega_A = 1$$

$$\omega_B = 1$$

$$\omega_C = 2$$

$$\Delta A = 0.2$$



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

w↓

Δ↑

Margin width ↑

$$\Delta A = 0.2$$

①

A (+)

$$\Delta c = 2$$

②

C (-)

B (+)
 $\Delta g = 5.8$

Δ_{\min}
 Δ_{\max}

Select

③

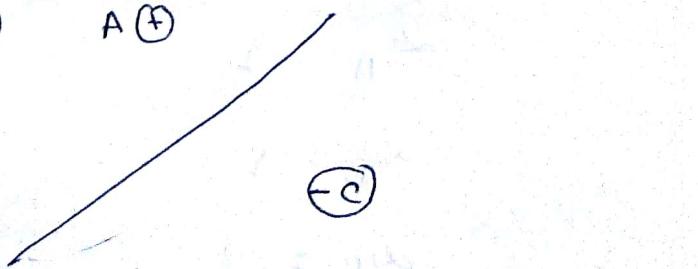
A (+)

B (+)

④

A (+)

C (-)



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)

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

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

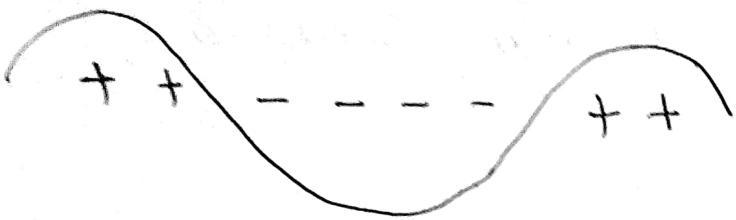
$$= \text{Sign} \left(\sum_i (\bar{w}_i y_i) K(\bar{x}_i \cdot \bar{x}) + b \right)$$

↳ 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



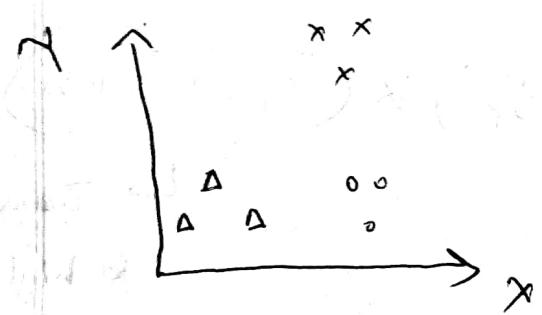
$$+ \circ + + \circ \circ + \circ + \circ +$$

— Polynomial

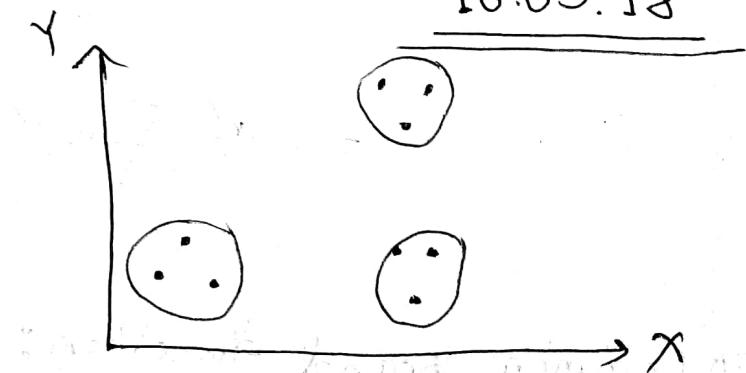
- Radio bias Kernel

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

$\longrightarrow x \longrightarrow$



Supervised



16.09.18
Unsupervised (not
labeled)

⇒ Supervised learning:

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

⇒ Unsupervised learning:

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

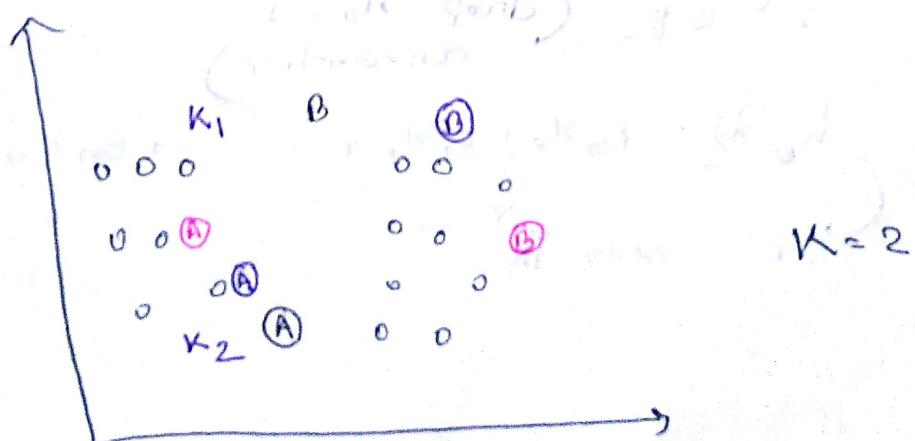
* Applications of clustering algorithms:

→ social network analysis

→ market segmentation

→ astronomical data analysis

* K-means Algorithm: $K = \text{No. of clusters.}$



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

* K-means:

→ iterative algo

→ has two steps

(i) Cluster assignment step
cluster

(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

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

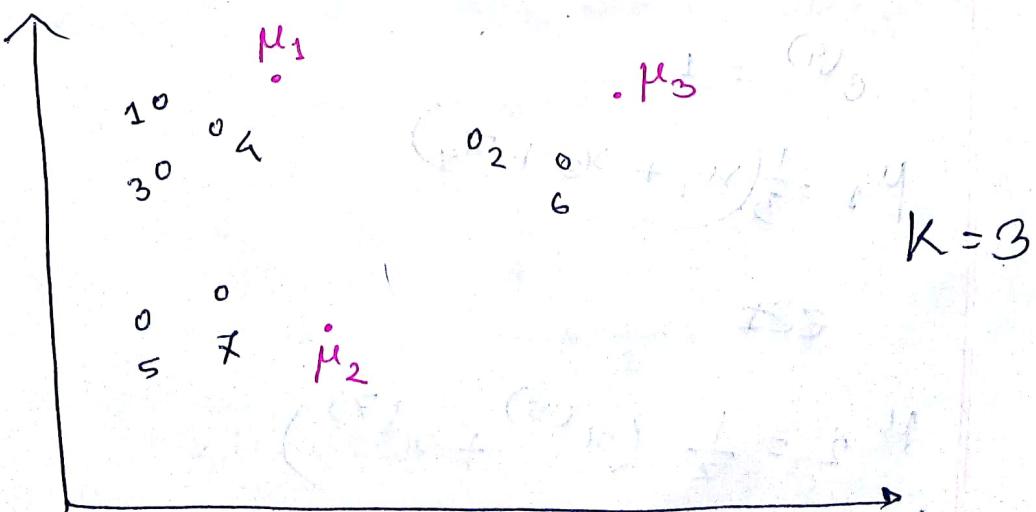
→ Repeat {

min J (-) w.r.t.
 for $i = 1$ to m
 $c^{(i)}$ = index (1 to K) of cluster centroid

min J (-) w.r.t.
 for $k = 1$ to K
 μ_k = average (mean) of points assigned
 to cluster k .

Cluster Assignment Step

Move centroid Step



$$X=3$$

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

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

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

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

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

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

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

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

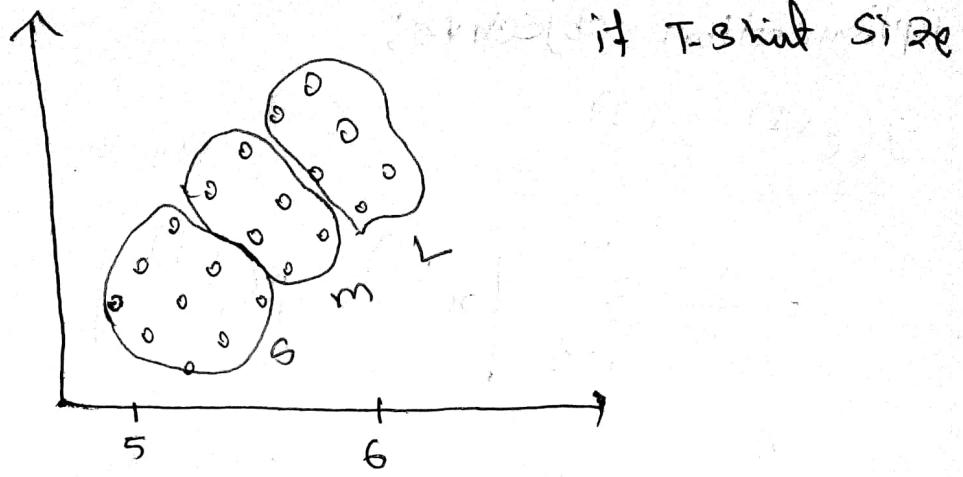
$$c^{(3)} = 1$$

$$c^{(4)} = 1$$

$$\mu_1 = \frac{1}{3}(x_1 + x_3 + x_4)$$

~~200~~

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

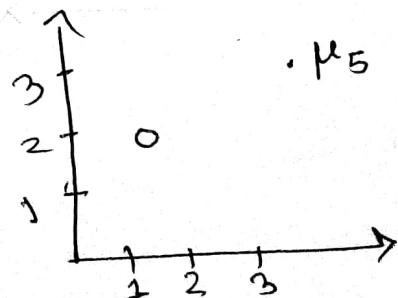


Optimization objective:

$c^{(i)}$ = 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,



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

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

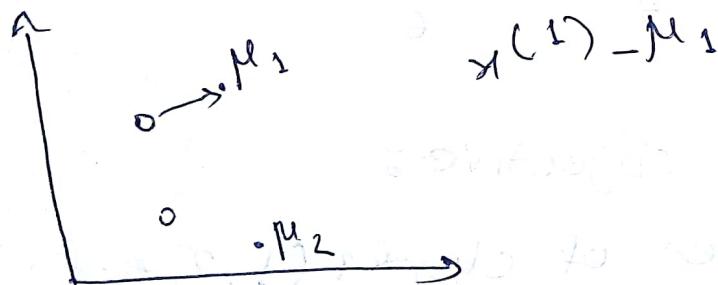
$$\therefore c^{(1)} = 5 \quad \mu_c^{(1)} = \mu_5 = (2, 3)$$

K

Optimization Objective:

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

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



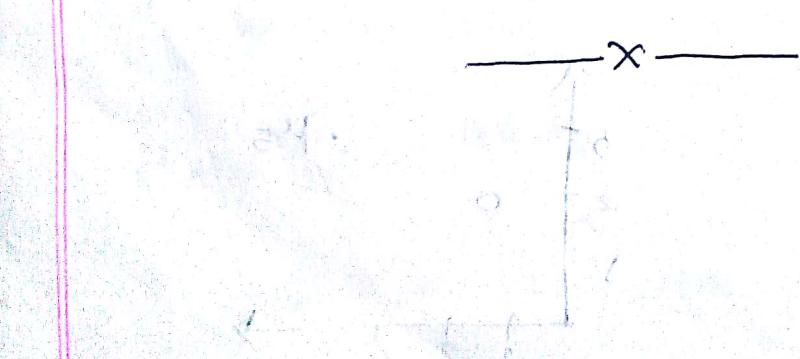
$$\min c^{(1)}, c^{(2)}, \dots, c^{(m)}$$

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

subject to $\mu_1, \mu_2, \dots, \mu_k$

values of signals to various methods

Surprise method with the signals



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:

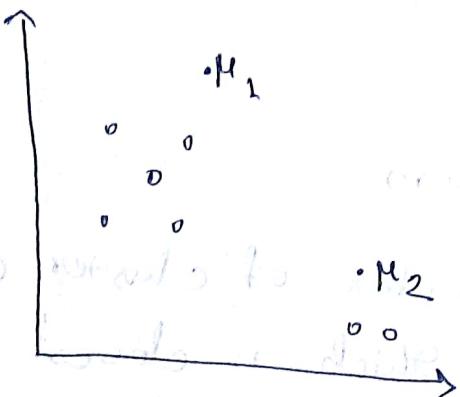
$$J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$
$$\min_{c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_k}$$

BT 09.81

ith index point $\rightarrow \mu_5$

blurred image $e^{(i)} = 5$

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



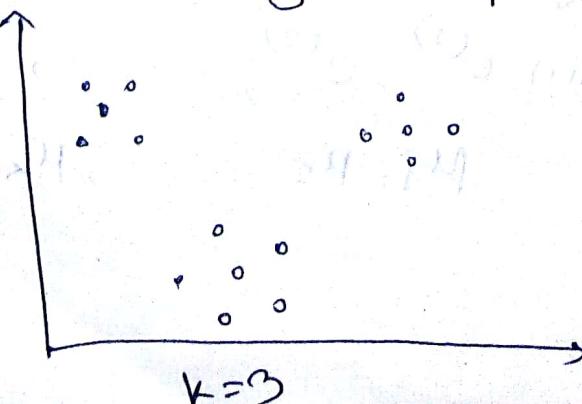
Important:

* 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)

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

}

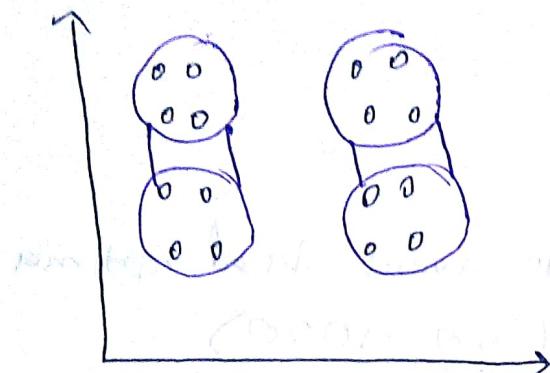
To pick the cluster that gives the lower cost value.

* (If $K = 2$ to 10) \uparrow this algo. works really well.

* If $K = 100$ + expected result can be get within few iterations.

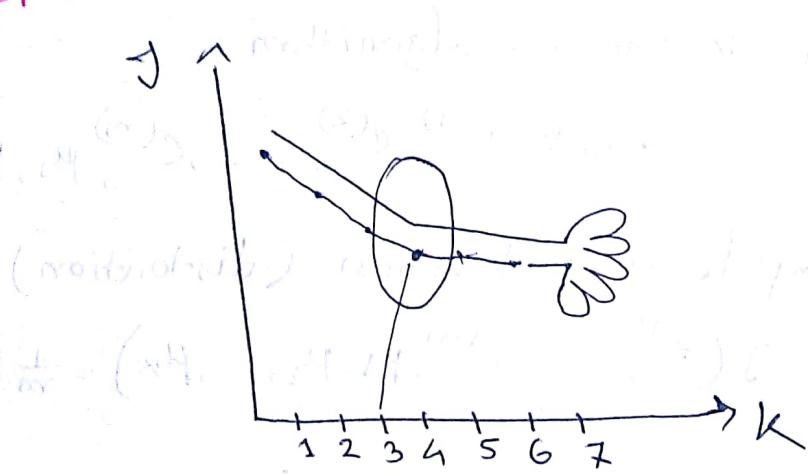
iterations. So, no. of iterations can be decreased.

How to choose value of K ?

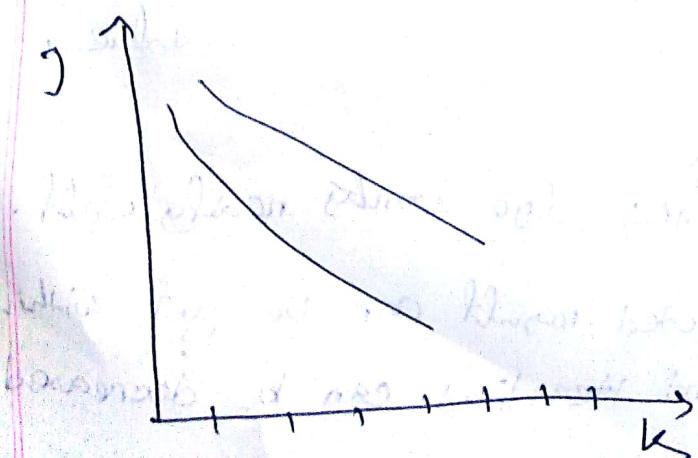


$$K = 4 / 2$$

* Elbow Method:



Select 3 as the value of K.



→ Elbow method

doesn't work
well in these
kinds of cases

* choosing K's value based on purpose.

