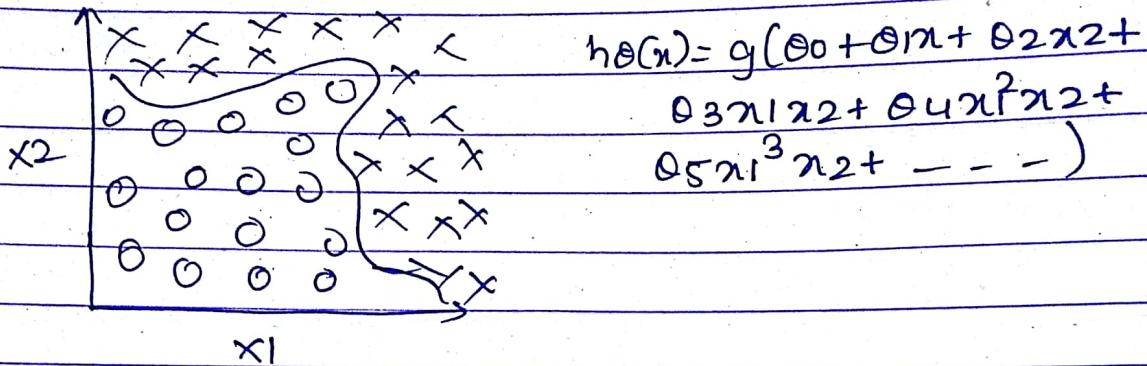


Unit - II

Supervised learning: Neural Networks

Considering a Non linear classification ex.



considering all quadratic terms

$$x_1^2, x_1 x_2, x_1 x_3, \dots, x_1 x_{100} \quad \left\{ \approx \frac{n^2}{2} \right.$$

we get such 5000 features

which is huge set of features and may over fit the data it is not a generalized equation
this increases implementing complexity

\therefore Including only subset of these features

$$\text{i.e. } x_1^2, x_2^2, x_3^2, x_4^2, \dots, x_{100}^2$$

No. of features are reduced i.e. 100

these features are very less and may not fit the data

considering 3rd order polynomial feature

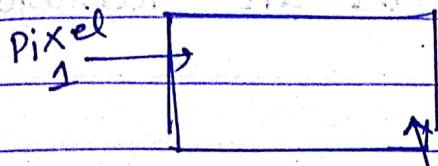
$$x_1 x_2 x_3, x_1^2 x_2, x_1 x_2 x_4, \dots$$

the features may range to 170,000

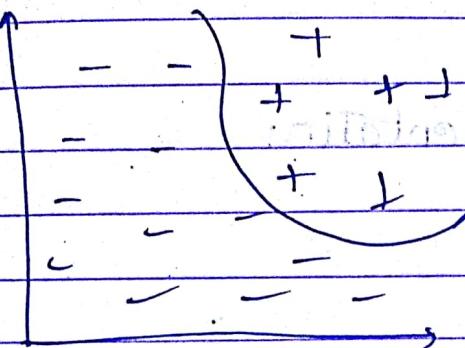
the features dramatically increases and additional features may cause problem and there may be issue of over fitting the hypothesis

reducing n reduces the features but most of the ML examples have large value of n

e.g.:



pixel 2



+ car

- Not car

image pixels $50 \times 50 = 2500$ pixels

~~$n=2500$~~ - gray scale

~~$n=7500$~~ if RGB

$n=$	pixel 1 intensity	pixel 2 intensity
"	2 "	"
"	3 "	"
"	!	"
"	2500	"

Quadratic features $(n \times n^2) \approx 3\text{million features}$

Neural Networks

- Algorithms that try to mimic the brain
- widely used in 80's and early 90's
- diminished in late 90's
- Recent resurgence: state of the art technique for many applications

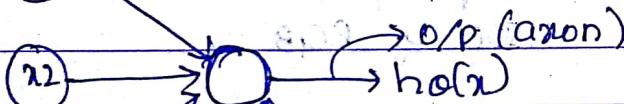
Neural networking exp (ppt)

Neural Networks Model Representation

Neuron in brain (ppt)

Neuron model: logistic unit

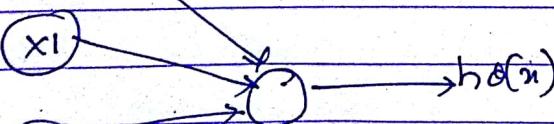
represents single neuron (comparing with brain neuron)



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

dendrites
"Input wires"

x_0 "bias unit"
 $x_0 = 1$



sigmoid (logistic) activation function

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

Neural
diff

x1

x2

x3

layer 1

"input layer"

$a_i^{(j)} =$

$e^{\theta^{(j)} x}$

Mathematica

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

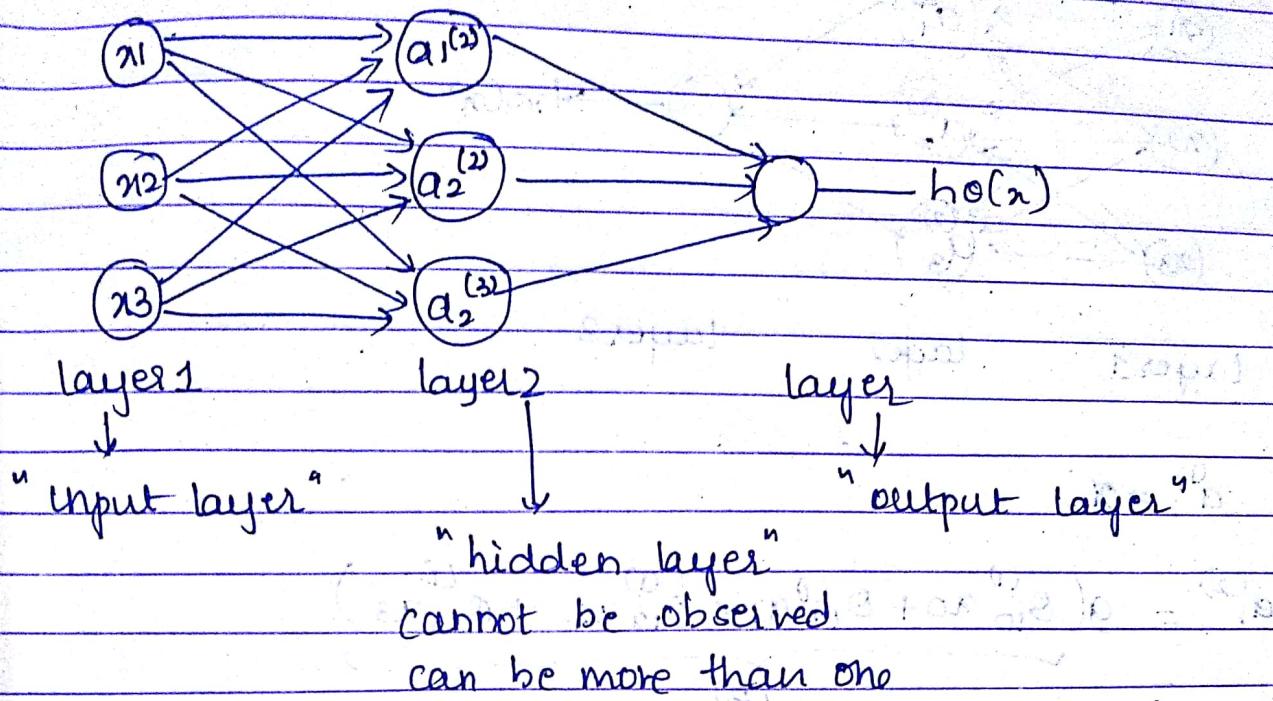
$a_2^{(2)} = g[$

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

$h\theta(x) =$

if n/10
Sj+1 un
dinner

Neural Net
different neurons strung together



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

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

Mathematical definition

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

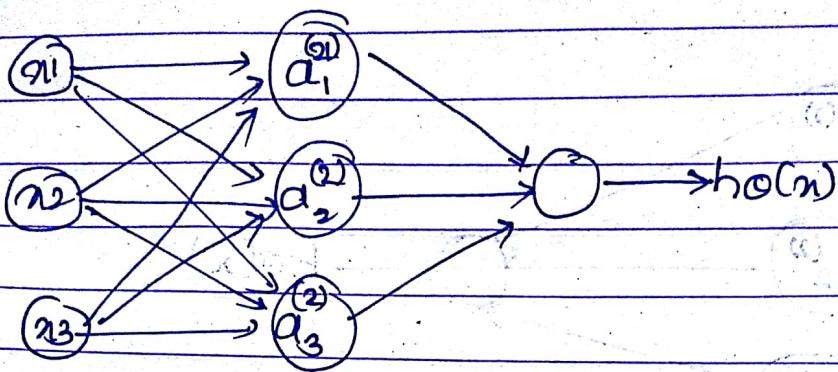
$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

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

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

if n/w has S_j units in layer j ,
S $_j+1$ units layer $j+1$ then $\theta^{(j)}$ will be of
dimension $(S_j+1) \times (S_j + 1)$

Forward propagation: Vectorized implementation



layer 1 layer 2 layer 3

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

$$a_1^{(2)} = g\left(\underbrace{z_1^{(2)}}_{= \theta_{10}x_0 + \theta_{11}x_1 + \theta_{12}x_2 + \theta_{13}x_3}\right)$$

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

$$a_2^{(2)} = g\left(\underbrace{z_2^{(2)}}_{= \theta_{20}x_0 + \theta_{21}x_1 + \theta_{22}x_2 + \theta_{23}x_3}\right)$$

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

$$a_3^{(2)} = g\left(\underbrace{z_3^{(2)}}_{= \theta_{30}x_0 + \theta_{31}x_1 + \theta_{32}x_2 + \theta_{33}x_3}\right)$$

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

$$h(\alpha(x)) = g\left(\alpha_0^{(2)} + \alpha_1^{(2)} + \alpha_{12}^{(2)} + \alpha_{13}^{(2)} \right)$$

$\underbrace{\hspace{10em}}$

$\alpha^{(3)}$
 z_1

$$\alpha = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix}$$

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

layer 2

$$z^{(2)} = \Theta^{(1)} \alpha. \quad \Rightarrow \quad z^{(2)} = \Theta^{(1)} \cdot \alpha^{(1)}$$

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

Kolmogorov

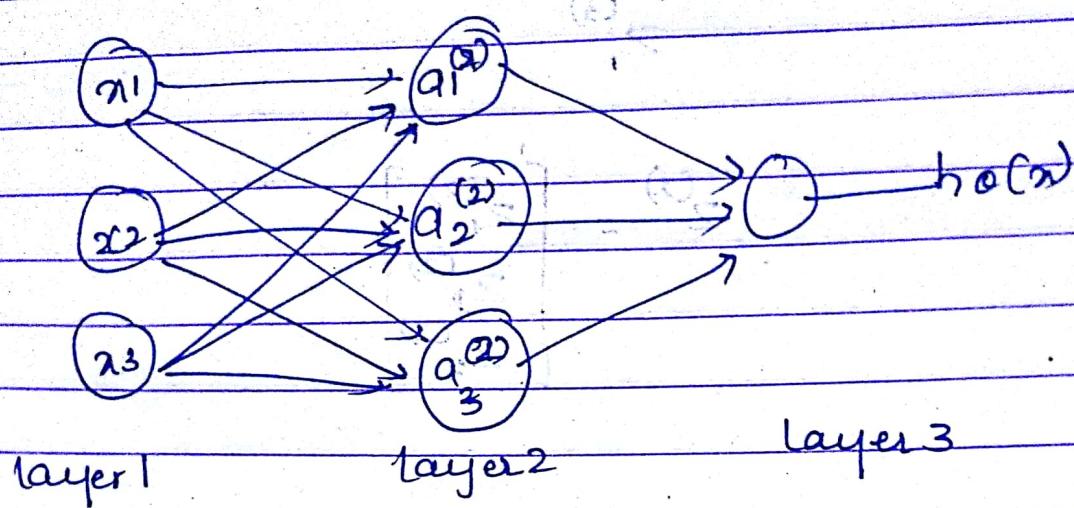
layer 3

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

$$\alpha^{(3)} = h(\alpha) = g(z^{(3)})$$

Forward propagation bcoz it is started off by activating the input unit then the hidden layer and the output layer.

Neural Networks learning its own features



considering only layer 2 and 3.

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

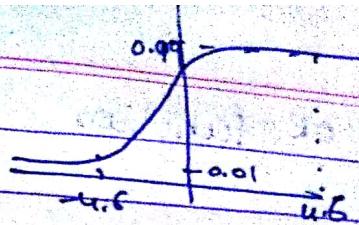
↳ logistic Regression

$$\begin{matrix} x_1, x_2, x_3. \\ a_1^{(2)} \quad a_2^{(2)} \quad a_3^{(2)}. \end{matrix}$$

Instead of taking quadratic equation or any other polynomial and restricting the learning of features.

The $\theta_j, \theta^{(j)}$ gives the flexibility of training the algorithm by itself.

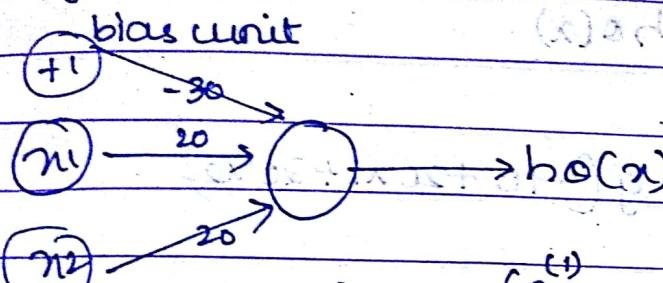
Example.



→ AND

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

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

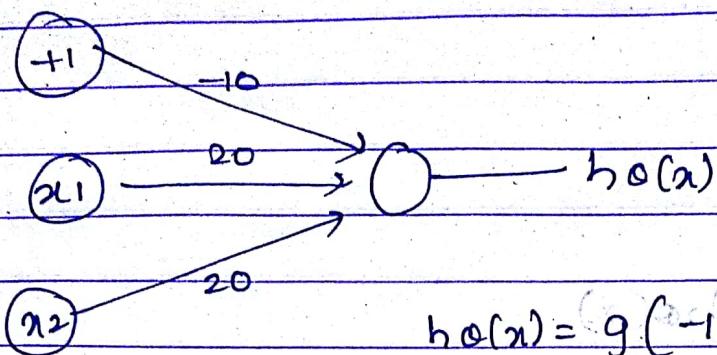


x_1	x_2	$h(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

$$h(x) \approx x_1 \text{ AND } x_2$$

FUNCTION $g(x)$ part

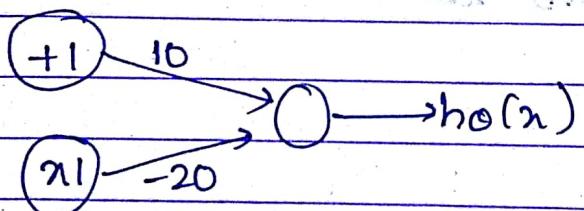
OR function



x_1	x_2	$h_{\theta}(x)$
0	0	$g(-10)$ ≈ 0
0	1	$g(10)$ ≈ 1
1	0	$g(10)$ ≈ 1
1	1	$g(30)$ ≈ 1

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

NOT function

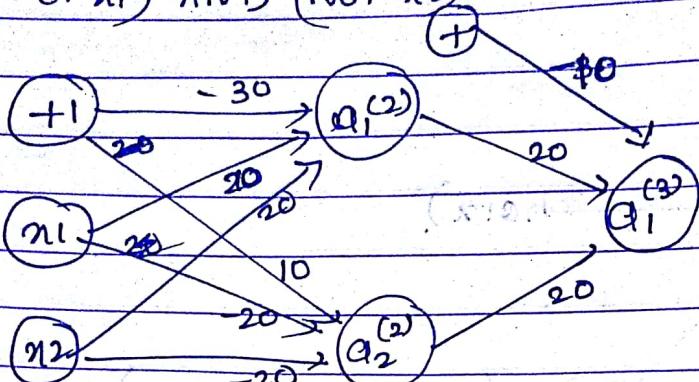


x_1	$h_{\theta}(x)$
0	$g(10)$ ≈ 1
1	$g(-10)$ ≈ 0

$h_{\theta}(x) \approx \text{NOT } x_1$

$\pi_1 \text{ XNOR } \pi_2$

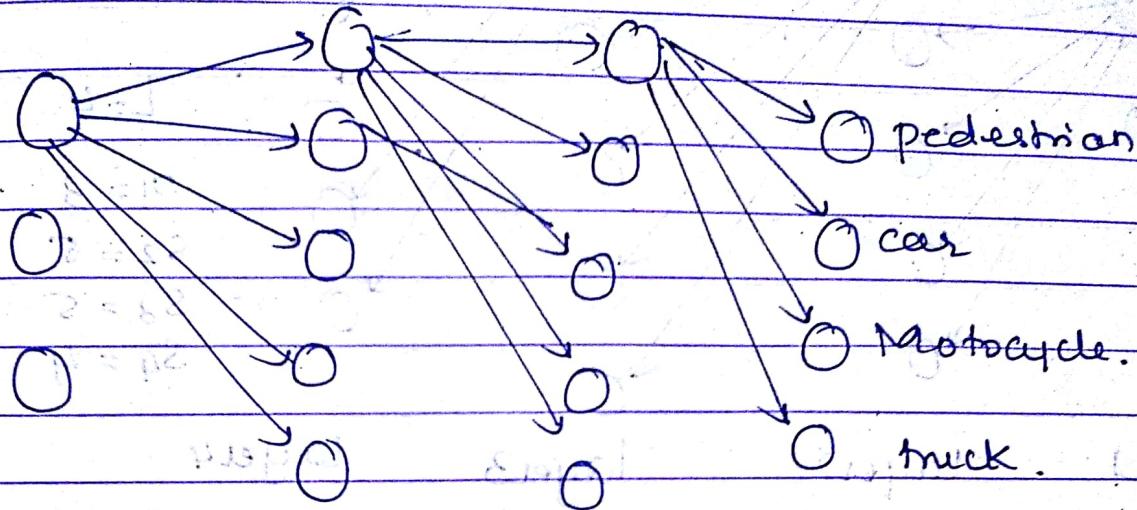
(NOT π_1) AND (NOT π_2)



π_1	π_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_0(\pi)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

Multi class Classification (One-vs-all)

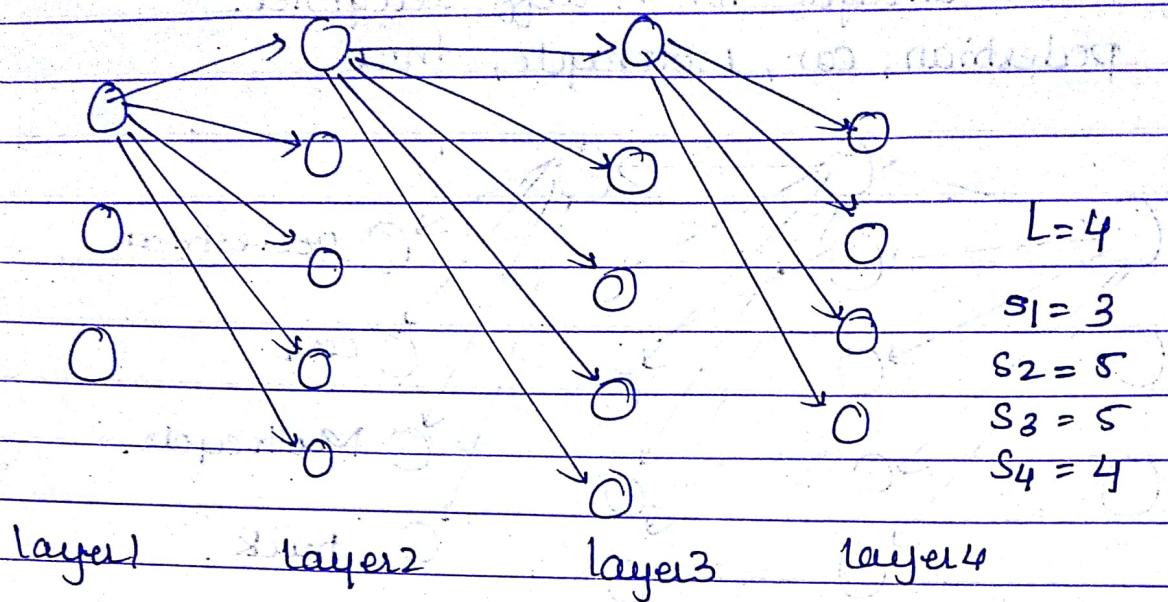
Classifying the images for 4 diff categories
pedestrian, car, motorcycle, truck.



$y^{(n)}$ is one of

	pedestrian	car	motorcycle	truck
pedestrian	1	0	0	0
car	0	1	0	0
motorcycle	0	0	1	0
truck	0	0	0	1

Representation



layer1 layer2 layer3 layer4

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

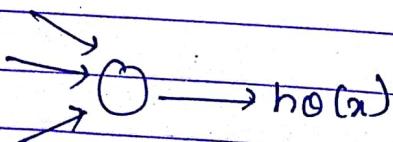
L = total No. of layers

s_L = no. of units in layer L

Binary Classification

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

1 output unit



$$s_L = s_1 = 1$$

k = also denotes no. of units in layer

Multi class classification (k classes)

if $K \geq 3$ we use one-vs-all method.

Cost function

generalized form of logistic regression.

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

Neural N/w :-

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

Need to compute

$$J(\theta)$$

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

Gradient computation

Back propagation algorithm

Intuition: $\delta_j^{(4)}$ = "error" of node j in layer 4

For each output unit (layer $L=4$)

$$\delta_j^{(4)} = a_j^{(4)} - y_j$$

$$(b_0^{(x)})_j \quad \delta^{(4)} = a^{(4)} - y$$

$$- \quad \delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} * g'(z^{(3)}) \rightarrow a^{(3)} * (1-a^{(3)})$$

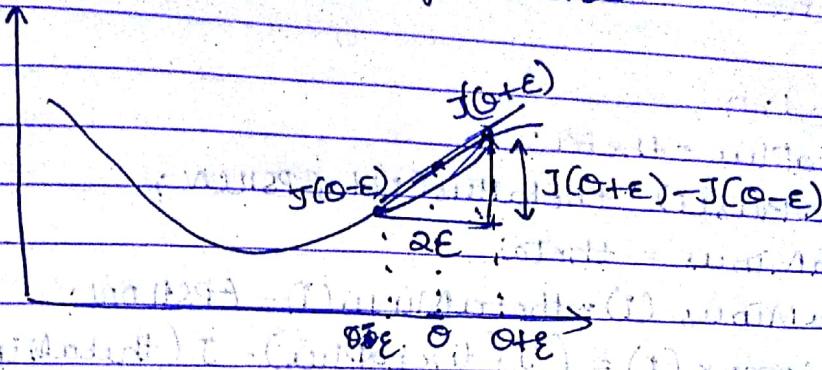
$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} * g'(z^{(2)}) \rightarrow a^{(2)} * (1-a^{(2)})$$

The name back propagation bcoz if s the
algorithm starts computing δ from output layer
through input layer.

ignoring λ if $\lambda=0$

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

Numerical estimations of gradients



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

$$\epsilon = 10^{-4}$$

ϵ value is very small so that the derivative term is approximately equal to the actual slope

Implement:

$$\text{gradApprox} = \frac{(J(\theta_1 + \epsilon) - J(\theta_1 - \epsilon))}{(2 * \epsilon)}$$

Parameter vector θ :

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

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

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

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

implement

for $i = 1:n$,

$\theta_{i+1} = \theta_i$;

$\theta_{i+1}(i) = \theta_i(i) + \epsilon$;

$\theta_{i-1} = \theta_i$;

$\theta_{i-1}(i) = \theta_i(i) - \epsilon$;

$\text{gradApprox}(i) = \frac{J(\theta_{i+1}) - J(\theta_{i-1})}{2\epsilon}$

end;

check $\text{gradApprox} \approx \text{DVec}$

Back propagation

if gradApprox is approx equal to DVec then
the algorithm implemented is correct.

Implementation Note

- Implement backprop to compute DVec

- Implement numerical gradient check to
compute gradApprox

- Make sure they give similar values

- Turn off gradient checking. Using backprop
code for learning

Important

BC sure to disable your gradient checking
code for before training your classifier. If you
run numerical gradient computation on every
iteration of gradient descent then your
code will be very slow

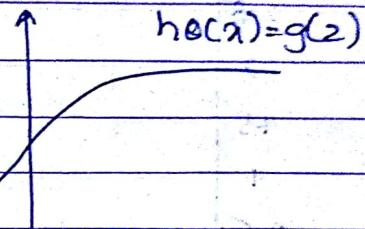
Autonomous vehicle (ex)

Support Vector Machines [SVM]

Objective optimization

In logistic regression

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



If $y=1$, then $h_0(x) \approx 1 \quad \theta^T x > 0$

If $y=0$ then $h_0(x) \approx 0 \quad \theta^T x < 0$

Cost function

alternative view

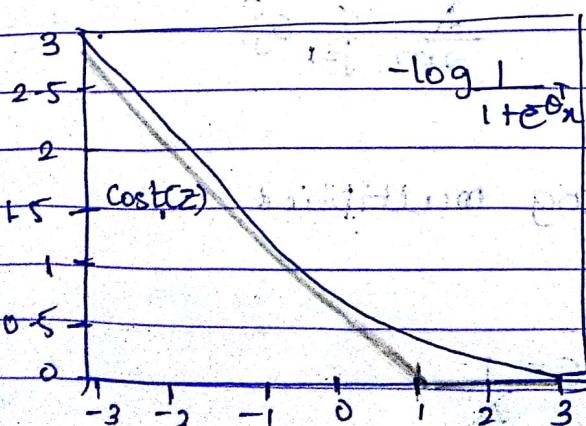
$$J(\theta) = -[y \log h_0(x) + (1-y) \log (1-h_0(x))]$$

↑ extra term contributed.

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

Considering two scenarios

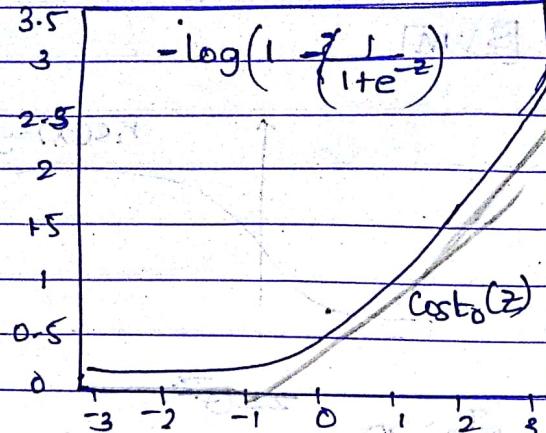
- i) $y=1 \quad (\theta^T x > 0)$



— cost function of SVM
— " " " of Logistic Regression

ii) if $y=0$

$$\theta^T x < 0$$



logistic regression

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

SVM:

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

m - constant value being multiplied

$$A + \lambda B$$

$$CA + \lambda B \quad i.e. C = 1$$

\therefore Removing m and adding c

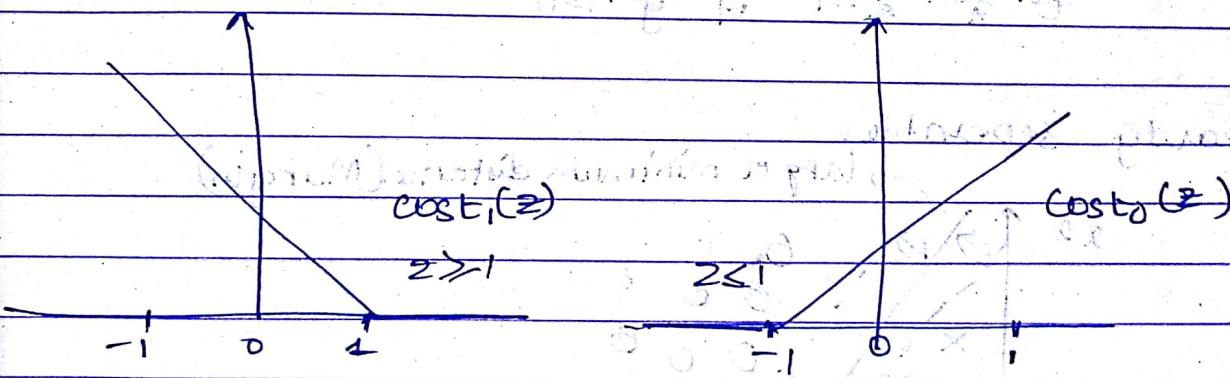
$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right]$$

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

hypothesis

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

Large Margin Intuition



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

if $y=0$ ($\theta^T x \leq -1$) not just < 0

~~svm decision boundary~~

~~consider very large C~~

$C = 100,000$

~~outcomes~~

~~softmax function~~

~~softmax~~

SVM Decision Boundary.

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

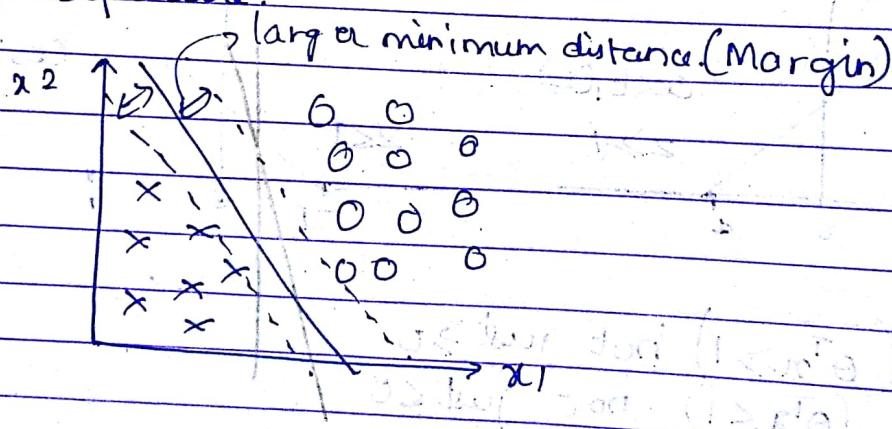
$$+ \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

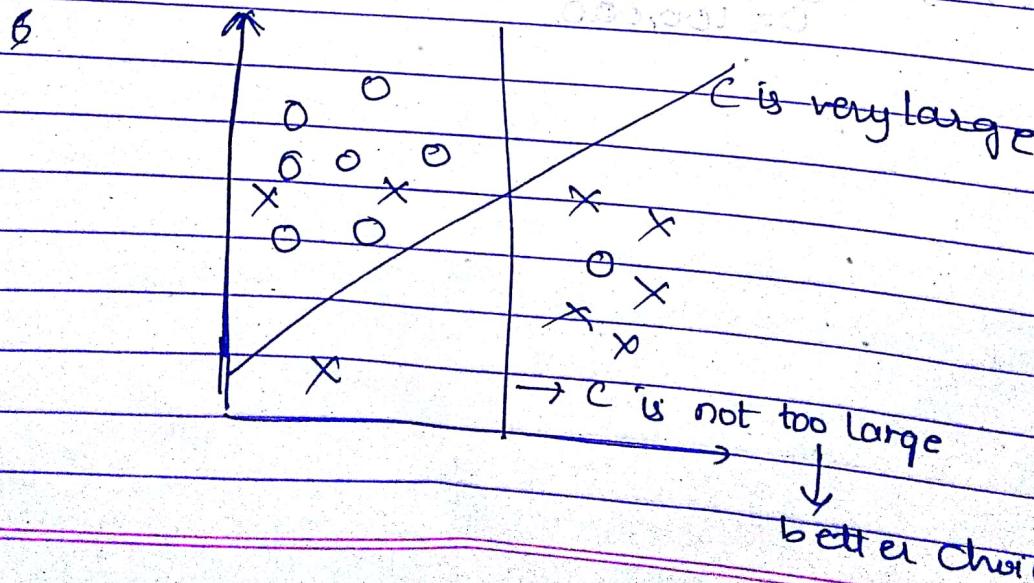
$$\text{ST. } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1$$

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

linearly separable.

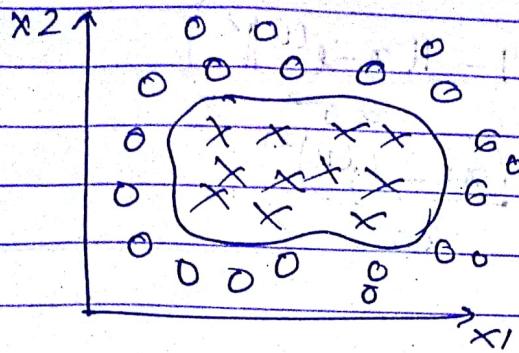


large margin classifier in presence of outliers



$\|w\|$ length of vector w

Kernels 1

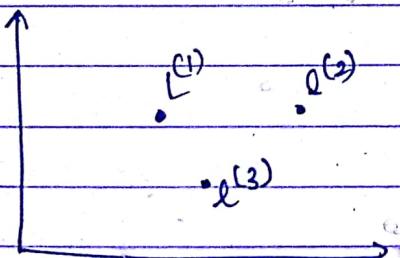


$$\text{predict } y=1 \text{ if } \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_1 x_2 + \alpha_4 x_1^2 + \alpha_5 x_2^2 + \dots \geq 0$$

$$h(x) = \begin{cases} 1 & \text{if quadratic} \\ 0 & \text{otherwise.} \end{cases}$$

$$\rightarrow \alpha_0 + \alpha_1 f_1 + \alpha_2 f_2 + \alpha_3 f_3 + \dots$$

Given x , compute new feature depending on proximity to landmarks $L^{(1)}, L^{(2)}, L^{(3)}$



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

$$f_2 = \text{similarity}(x, L^{(2)}) = \exp\left(1 - \frac{\|x - L^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \text{similarity}(x, L^{(3)}) = \exp\left(1 - \frac{\|x - L^{(3)}\|^2}{2\sigma^2}\right)$$

kernel

(Gaussian kernel)

Kernels and similarity

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

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

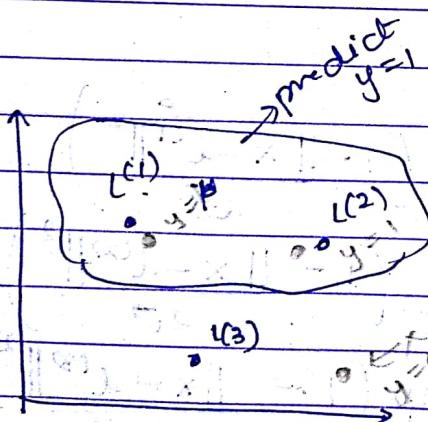
If $x \approx l^{(1)}$:

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

If x_0 is far from $l^{(1)}$:

$$f = \exp\left(-\frac{(\text{large } x_0)^2}{2\sigma^2}\right) \approx 0$$

Ex:



predict 1 when

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

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

$$f_1 \geq 1 \quad f_2 \geq 0 \quad f_3 \geq 0$$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 = -0.5 + 1 = 0.5 \geq 0$$

predict $y=1$

$$f_1 \neq f_2 \neq f_3 \neq 0$$

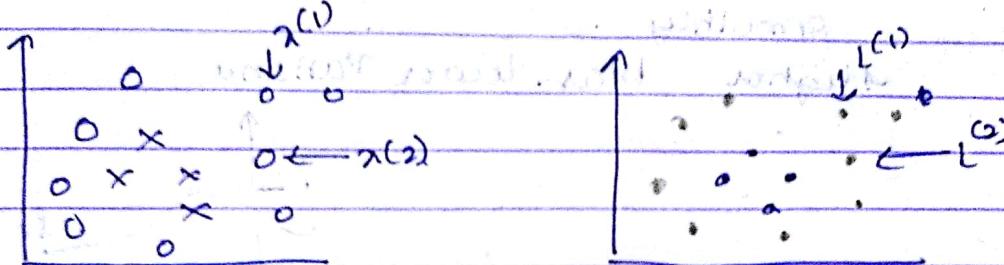
$$\alpha_0 + \alpha_1 f_1 + \alpha_2 f_2 + \alpha_3 f_3 = -0.5 \leq 0$$

predict: $y=0$

SVM with kernels II

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

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



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

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

⋮

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

$$f_1^{(p)} = \text{sim}(x^{(i)}, L^{(1)})$$

$$f_2^{(p)} = \text{sim}(x^{(i)}, L^{(2)})$$

$$\vdots \rightarrow f_i^{(p)} = \text{sim}(x^{(i)}, L^{(i)}) = c,$$

$$f_m^{(p)} = \text{sim}(x^{(i)}, L^{(m)})$$

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$$

* half

SVM parameters

outfitting
 $C = \lambda$

hyp large C : lower bias, high variance

(small)

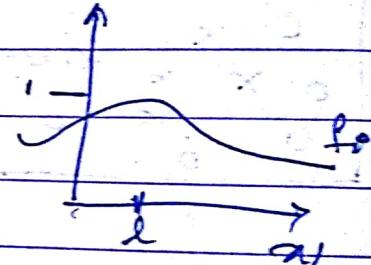
hyp small C : higher bias, low variance

(large)

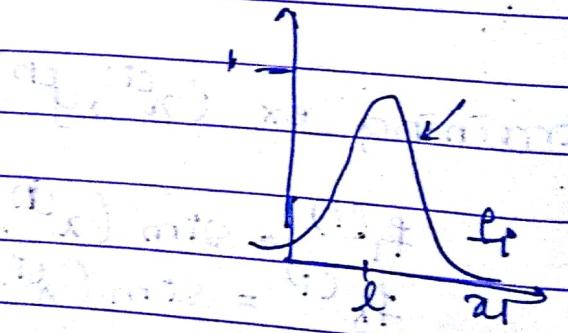
under fitting

σ^2 : large σ^2 , features fit vary more smoothly

higher bias, lower variance

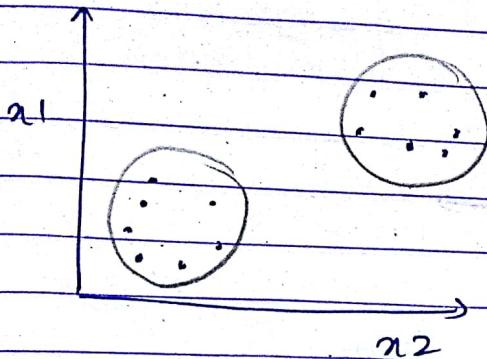


small σ^2 : features fit vary less, smoothly
lower bias high variance



Unsupervised learning

Clustering



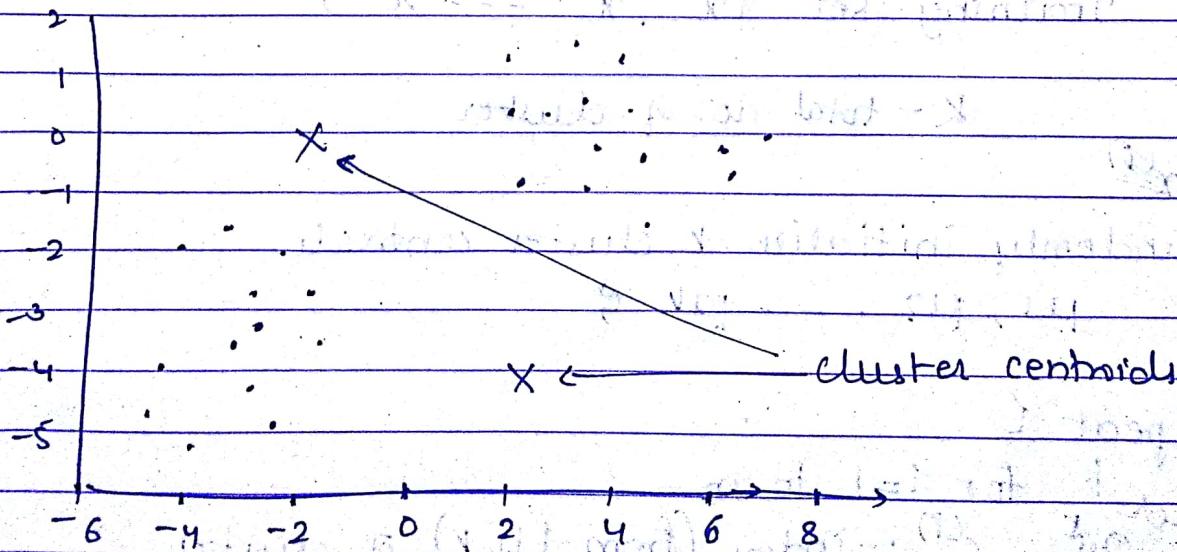
Clustering
algorithm.

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

Applications

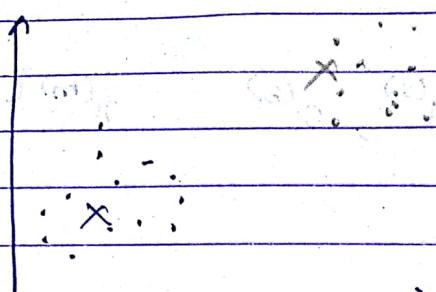
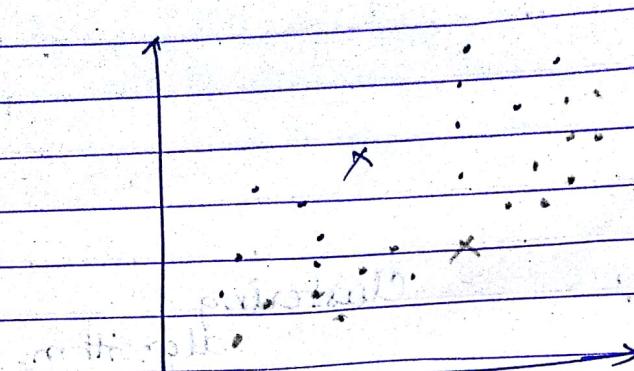
- 1) Market segmentation
- 2) Social network analysis
- 3) Organize computing clusters
- 4) Astronomical data analysis

K-means Algorithm (Iterative algorithm)



- 1) Cluster assignment set
- 2) Move centroid set

K-means



multimodal. local minima
and local maxima

optim

$C^{(p)}$

μ_k

J_{dc}

K-means algorithm

Input

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

K - total no. of cluster

* randomly initialize k cluster centroids
 $\mu_1, \mu_2, \dots, \mu_k$

repeat {

cluster for $i=1$ to n
assign $x^{(i)}$ to cluster $c^{(p)}$: index (from 1 to k) of cluster
set centroid closest to $x^{(i)}$
minimizing $J(\dots)$ wrt $c^{(1)}, c^{(2)}, \dots, c^{(m)}$

more for $k=1$ to K

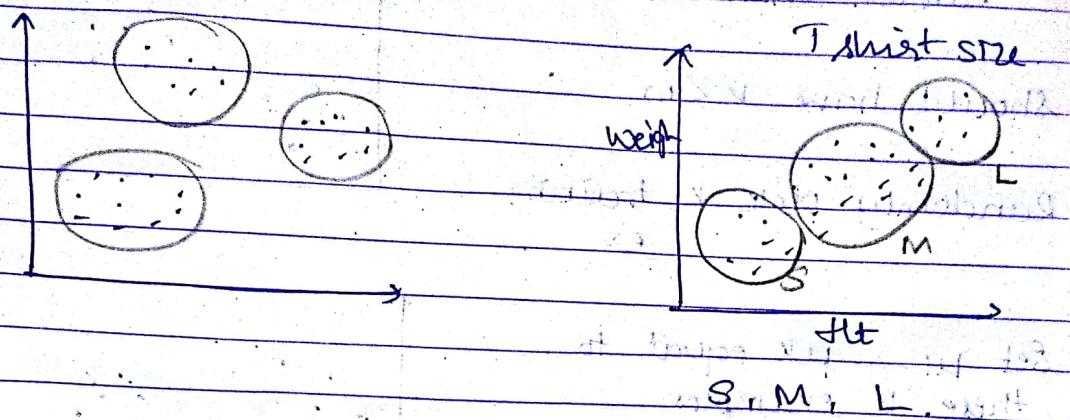
centroid $\mu_k := \text{avg}(\text{mean})$ of points assigned to cluster k
minimizing $J(\dots)$ wrt μ_1, \dots, μ_k

Optimi

J
(i) min

μ_1, \dots

K-means for non separated clusters.



Optimization Objective (for cost function)

$c^{(i)}$ = index of cluster

μ_k = cluster centroid k

$d_c(i)$ = cluster centroid of cluster to which

ex $x^{(i)}$ has been assigned.

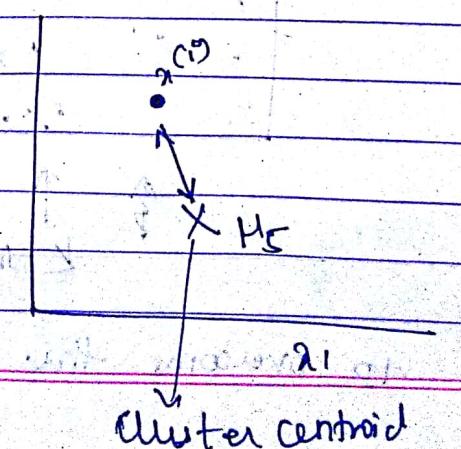
$$x^{(p)} \rightarrow 5 \quad c^{(i)} = 5 \quad \mu_c^{(i)} = \mu_5$$

Optimization objective:

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

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

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



** random initialization

should have $k < m$

Randomly pick k training
ex

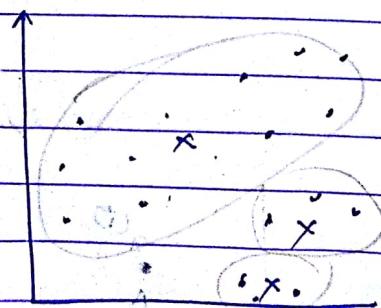
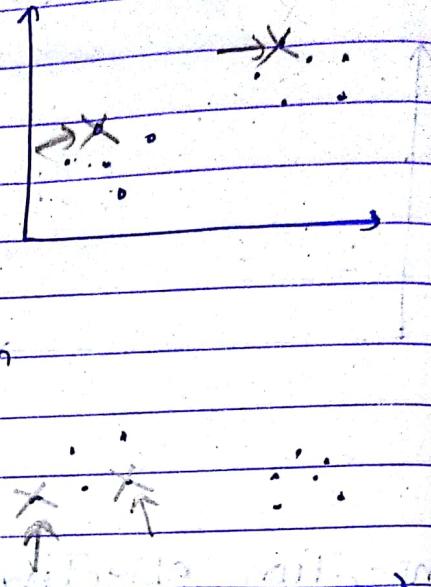
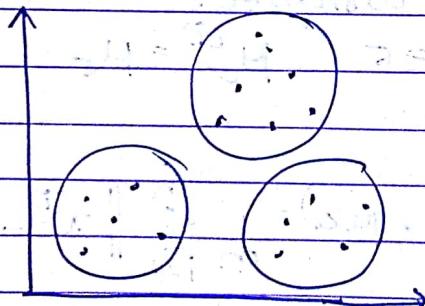
Set μ_1, \dots, μ_k equal to
these k examples

$k=2$

$$\mu_1 = x^{(i)}$$

$$\mu_2 = x^{(j)}$$

local optima



Best possible

cost function

k -means cluster is stuck @
a local minima
to overcome this ~~choose~~ initialise k many
times

Random initialization

For $i=1 \dots 100$

{ Randomly initialize k-means.

Run k-means,

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

$\mu^{(1)}, \dots, \mu^{(k)}$

Compute cost function (distortion)

$J(c^{(1)}, \dots, c^{(m)}, \mu^{(1)}, \dots, \mu^{(k)})$

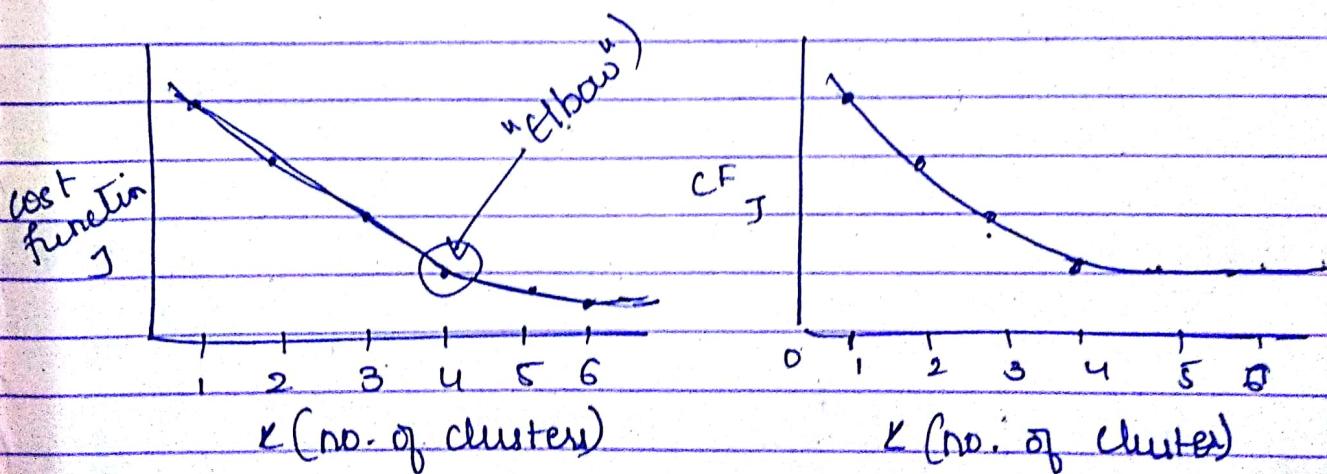
}

Pick clustering that gives lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu^{(1)}, \dots, \mu^{(m)})$

Choosing the number of clusters

→ Manually (Visualization)

Elbow Method:

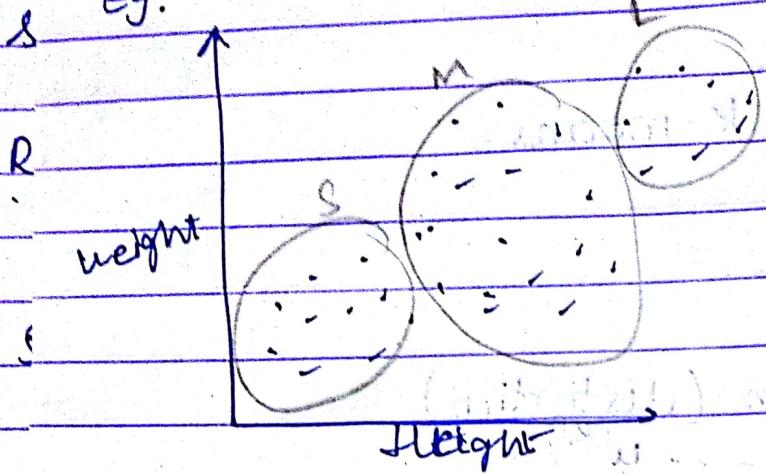


using 4 cluster
may be right
distortion goes down
at this point

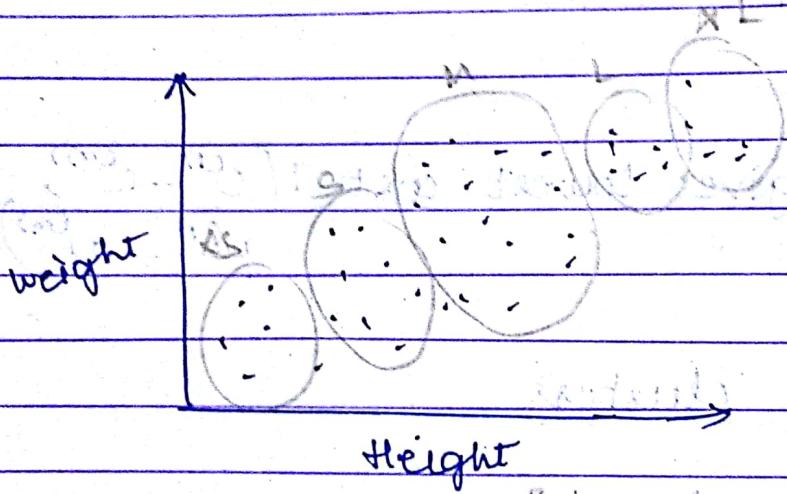
theres no clear
elbow there no
clear answer

* Metric:
purpose for running k-means

Eg:



$k=3$



$k=5$