

1. Supervised learning  
 2. Unsupervised learning

- Intro to

Andrew Ng

Supervised learning: learning from examples

- regression (house size - price problem)

continuous data of features is given. New data is placed in accordance with previous one.

Classification (cancer problem)

features are classified into 2 or more set outputs.

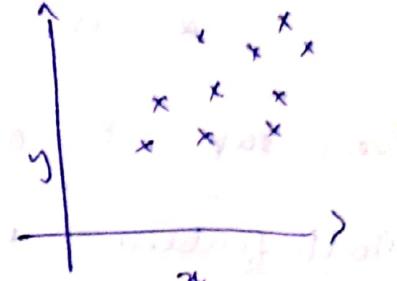
Unsupervised learning (clustering)

output of feature values are not specified. Algorithm divides values based on clusters.

### Linear Regression

training data ( $x$ )

output ( $y$ )



$x^{(i)}$  → first  $i^{\text{th}}$  training data

$y^{(i)}$  →  $i^{\text{th}}$  output data

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

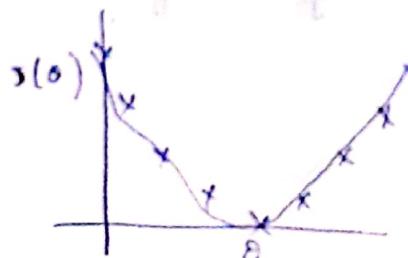
$\theta_0, \theta_1$  → parameters of model.

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

minimise error

(Squar error cost fn)

Start with some  $\theta$  value and plot  $J(\theta)$  vs  $\theta$



choose  $\theta$  for which  $J(\theta)$  is minimum

for more ~~parameters~~, contour plot one more.

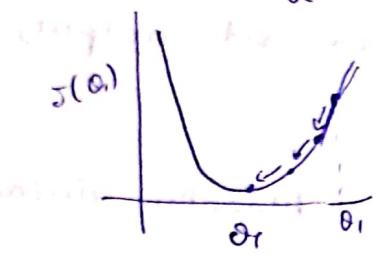
The values of  $\theta$ 's for which contour is at lower are the values to be chosen.

gradient descent: (Batch gradient descent)

to find  $\theta$  value

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

$\alpha \rightarrow$  learning rate. (determines number of steps to take)  
choose a  $\alpha$ ,



derivative term is given by tangent positive slope.  
So next value of  $\theta$  is lower.  
This will keep happening  
till  $\theta$  is min.

If  $\alpha$  is small, gradient descent is slow  
 $\alpha$  is large, gradient descent can overshoot

Every step in  $\theta$  OGD involves all training sets.

### Multi feature Linear Regression

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

let  $x_0 = 1$ .

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

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

$\theta$ -parameter vector

$x$ -feature vector.

Same gradient descent algorithm for them

feature scaling:

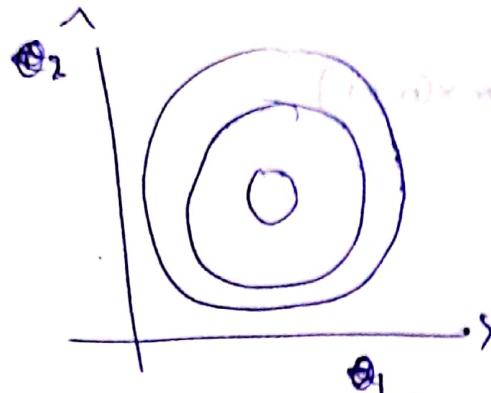
If our feature is more has larger range than others, then gradient descent will find it very slow.

$$x_1 \in (0 - 2000)$$

$$\text{let } x_1 = \frac{\text{value}}{2000}$$

$$x_2 \in (-5, 5)$$

$$x_2 = \frac{\text{value}}{5}$$



more features  $\rightarrow$  slow  
more  $f(x) = 0$

$$f(x^T)(x^T x) = 0$$

optimum value

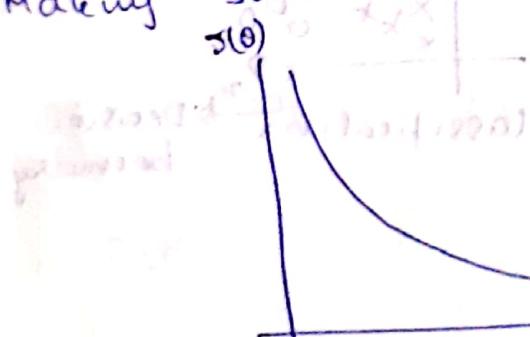
$$\rightarrow x \in 1.$$

Mean normalisation:

Replace  $x_i$  with  $x_i - \bar{x}_i$  to make features approximately zero mean.

$$x = \frac{x_i - \text{avg value}}{\text{range}}$$

Making sure  $a > 0$  works properly.



min value achieved.

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

Polynomial Regression:

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

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

Normal Equation method.

$$X = [x_0 \ x_1 \ x_2 \ \dots \ x_n] \ m \times (n+1)$$

$$y = [y] \ m \times 1$$

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

no need to choose  $\alpha$

no need to iterate

slow if  $n$  (no of feature) is large

$X^T X$  can be non-invertible when

- redundant feature

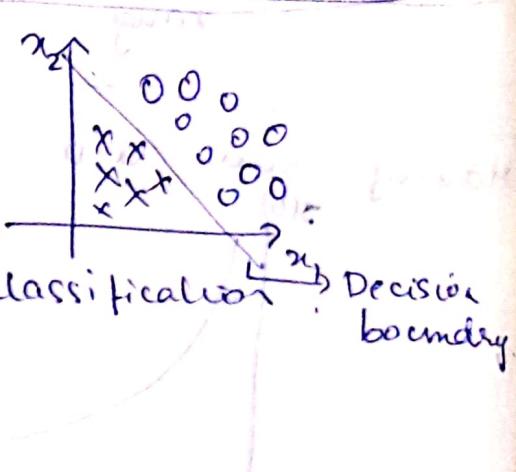
- lot of feature ( $m \leq n$ )

Classification:

Logistic Regression:

output  $y \in \{0, 1\}$  binary classification

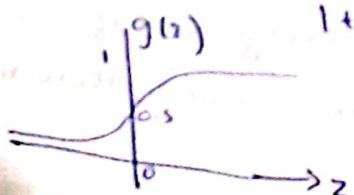
$$\text{OR } h_0(x) < 1$$



$$h_0(x) = g(\theta^T x), \quad g(\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n)$$

$$g(z) = \frac{1}{1 + e^{-z}} \rightarrow \text{logistic function}$$

Sigmoid function.

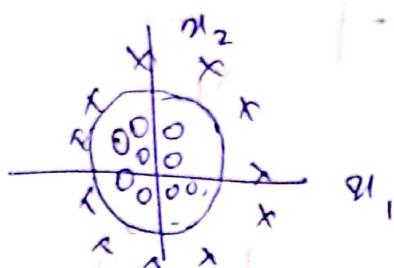


$$h_0(x) = P(y=1 | x, \theta)$$

probability that  $y=1$ , given  $x$ , parameterised by  $\theta$ .

$$\text{e.g. } h_0(x) = g(z) > 0.5 \text{ when } \theta^T x > 0$$

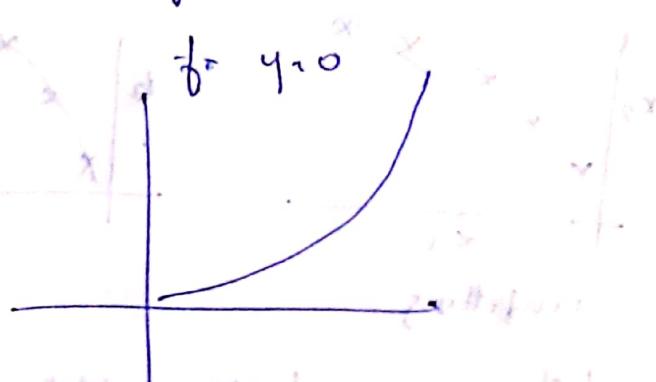
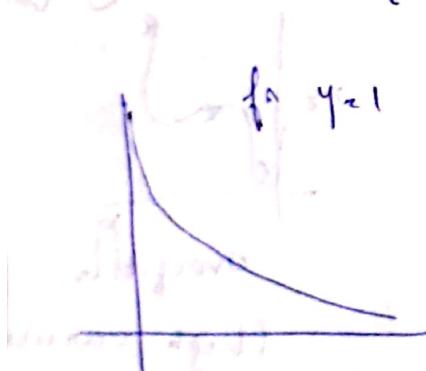
you linear decision boundary



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

cost function:  $-\log h_0(x)$  if  $y=1$   
 $-\log(1-h_0(x))$  if  $y=0$

$$\text{cost}(h_0(x), y) = \begin{cases} -\log h_0(x) & y=1 \\ -\log(1-h_0(x)) & y=0 \end{cases}$$



$$\therefore \text{cost}(h_0(x), y) = -y \log h_0(x) - (1-y) \log(1-h_0(x))$$

$$\therefore J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y_i \log h_0(x_i) + (1-y_i) \log(1-h_0(x_i))]$$

$$\text{with } \theta_j = \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \rightarrow \text{UND}$$

use various optimisation algorithms such as Conjugate gradient, BFGS, LBFGS.

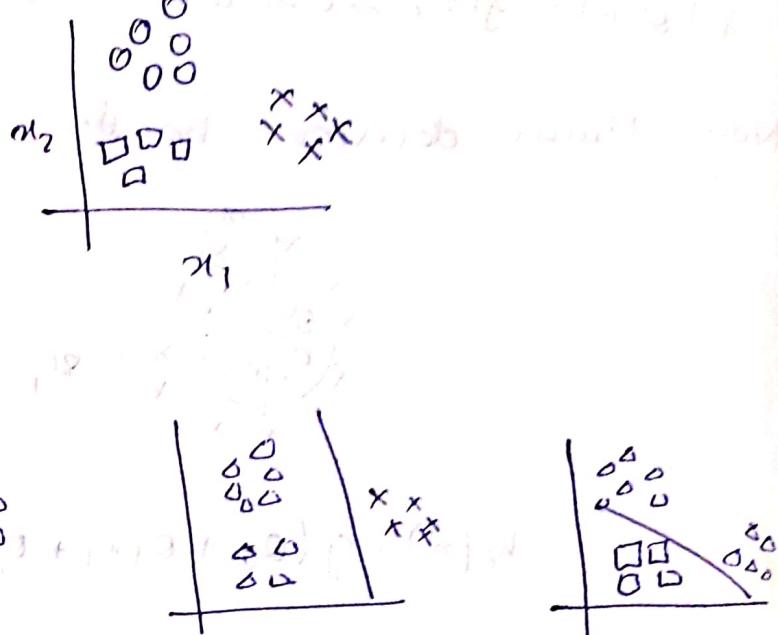
### Multi Class Classification

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

Class 1: 0

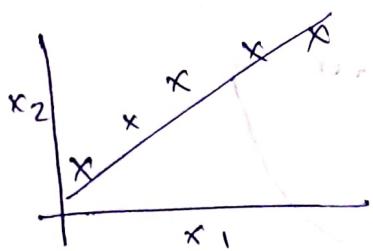
Class 2: □

Class 3: x

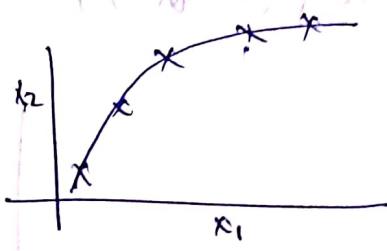


similarly do for class 2, 3.

### Overfitting



unfitting



overfitting  
(high variance)

overfitting can happen even though we use too

many features. lost function But it may fail to categorise new data.

### Regularisation

we can make some parameters  $\approx 0$ .

(penalising parameters)

Small values of  $\theta_0, \theta_1$  - less prone to overfitting.  
- Simple hypothesis.

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

helps in achieving this.

in linear regression.

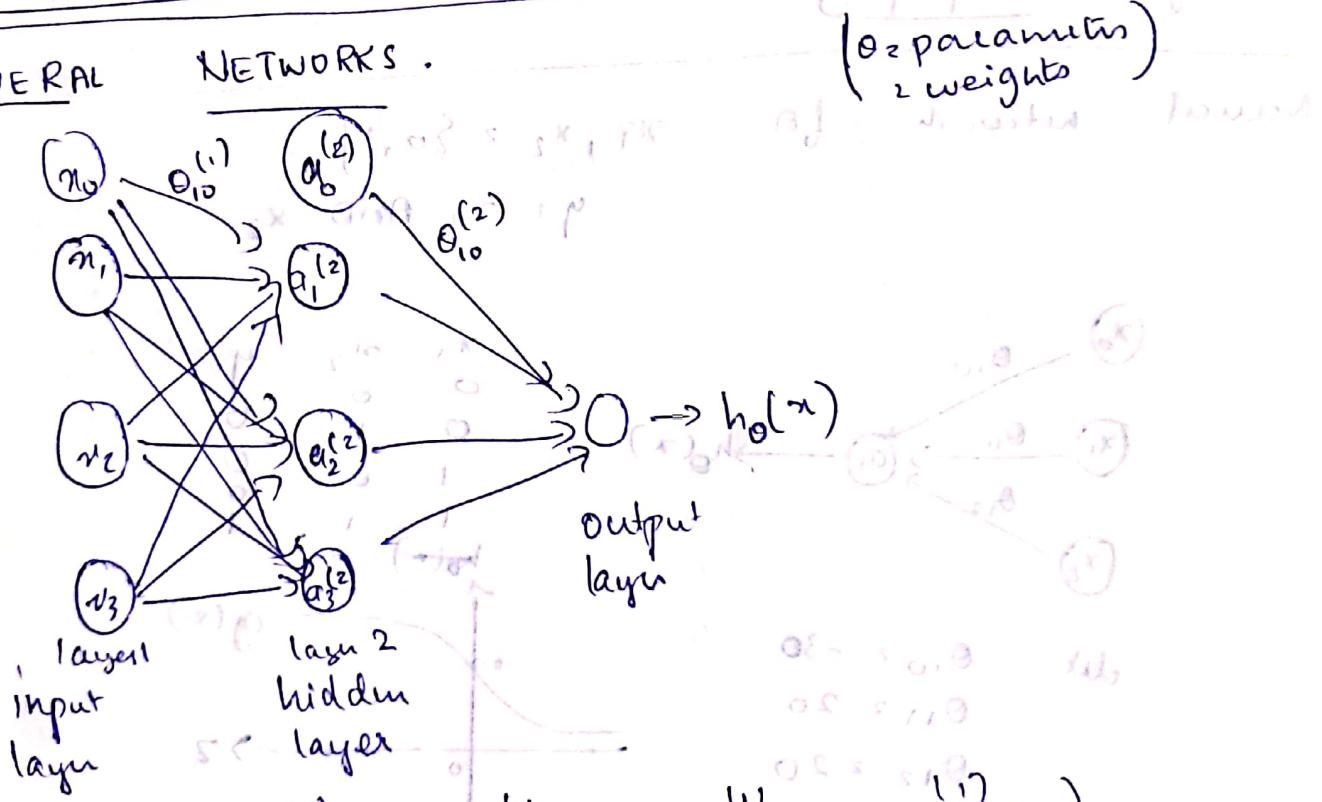
$$\theta_j = \theta_j \left(1 - \frac{\alpha x_j}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Shrinking parameter -

for logistic regression -

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

## GENERAL NETWORKS.



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

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

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

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

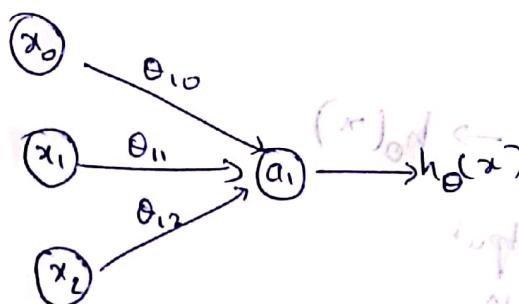
$$\begin{aligned}
 & \text{let } a_1^{(2)} = g(z_1^{(2)}) \\
 & a_2^{(2)} = g(z_2^{(2)}) \\
 & a_3^{(2)} = g(z_3^{(2)}) \\
 & z^{(2)} = \begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{pmatrix} \\
 & z^{(2)} = \theta^{(1)}x \\
 & a^{(2)} = g(z^{(2)}) \\
 & z^{(3)} = \theta^{(2)}a^{(2)} \\
 & h_{\theta}(x) = a^{(3)} = g(z^{(3)})
 \end{aligned}$$

forward propagation

Neural network for

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

$$y = a_1 \text{ AND } x_2$$

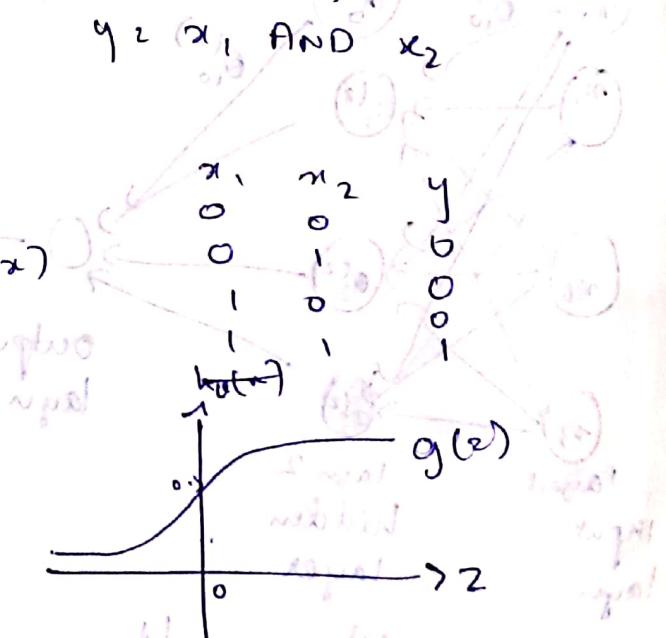


$$\text{let } \theta_{10} = -30$$

$$\theta_{11} = 20$$

$$\theta_{12} = 20$$

$$\begin{aligned}
 & x_1 \quad x_2 \quad h_{\theta}(x) \\
 & 0 \quad 0 \quad g(-30) \approx 0 \\
 & 0 \quad 1 \quad g(-10) \approx 0 \\
 & 1 \quad 0 \quad g(-10) \approx 0 \\
 & 1 \quad 1 \quad g(10) \approx 1
 \end{aligned}$$

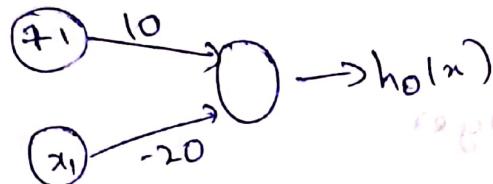


at network for OR,

$$\theta_{10} = 10 \quad \theta_{12} = 20 \quad \theta_{13} = 20$$

$x_1$	$x_2$	$h_{0(n)}$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$g(10) \approx 1$
1	1	$g(30) \approx 1$

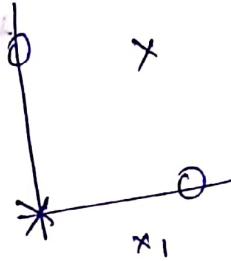
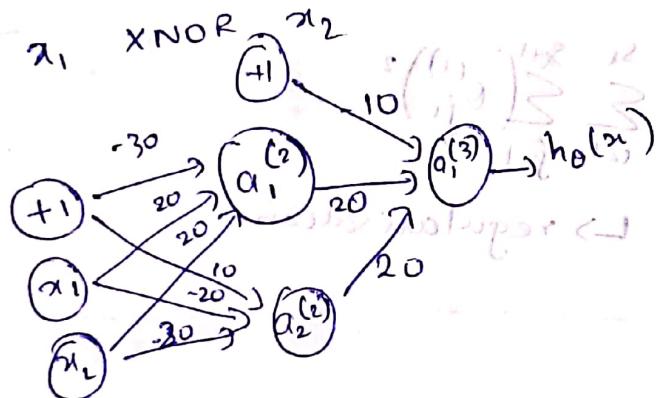
Neural network for NOT.



$x_1$	$h_{0(n)}$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

$$h_{0(n)} = g(10 - 20x_1)$$

(Utilizing AND of NOT)  $\Rightarrow$   $h_{0(n)} = g((x_1 \cdot \bar{x}_1) + 10)$



Inputs with performance

$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_{0(n)}$
0	0	0	0	0
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

$$(0,0) \oplus (0,0) = 0$$

$$(0,0) \oplus (0,1) = 0$$

$$(0,1) \oplus (0,0) = 0$$

$$(0,1) \oplus (1,1) = 1$$

Multiclass classification.  
No of outputs of Neural net should be no of classes.

(i),

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

different outputs

Cost function for Neural Network.

$L \rightarrow$  no of layers

$s_L \rightarrow$  no of neurons in a layer

for binary classification,

$s_L = 1$ ,

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

$\hookrightarrow$  regularisation

To minimise cost function.

Backpropagation algorithm

for a training example,

do forward propagation

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

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

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

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

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

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

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

$\delta_l$  = error of node  $j$  in layer  $l$ .

$$\delta_j^{(l)} = o_j^{(l)} - y_i \cdot h_0^{(x)} \quad \begin{array}{l} y_i \rightarrow \text{value of } y \text{ in} \\ \text{training set} \end{array}$$

$$\delta^{(1)} = o^{(1)} - y \quad \begin{array}{l} \text{what was output by} \\ \text{learning theory} \end{array}$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(1)} \cdot g'(z^{(2)}) \quad \begin{array}{l} \text{what was output by} \\ \text{learning theory} \end{array}$$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(2)} \cdot g'(z^{(3)}) \quad \begin{array}{l} \text{what was output by} \\ \text{learning theory} \end{array}$$

there is no  $\delta^{(1)}$ .  
 weight pattern  
 weight pattern  
 weight pattern

set  $\Delta_{ij}^{(1)} = 0$ .

$$\Delta_{ij}^{(1)} = \Delta_{ij}^{(1)} + o_j^{(1)} \delta_i^{(2)} \quad \begin{array}{l} \text{for accumulating } \delta \text{ terms.} \\ \text{if } j \neq 0 \text{ then } \Delta_{ij}^{(1)} = 0 \end{array}$$

$$\Delta_{ij}^{(2)} = \Delta_{ij}^{(2)} + \delta^{(2)} (o^{(2)})^T \quad \begin{array}{l} \text{if } j \neq 0 \text{ then } \Delta_{ij}^{(2)} = 0 \end{array}$$

$$D_{ij}^{(1)} = \frac{1}{m} \Delta_{ij}^{(1)} + \lambda \theta_{ij}^{(1)} \quad \begin{array}{l} \text{if } j \neq 0 \text{ then } D_{ij}^{(1)} = 0 \end{array}$$

$$D_{ij}^{(2)} = \frac{1}{m} \Delta_{ij}^{(2)} \quad \begin{array}{l} \text{if } j \neq 0 \text{ then } D_{ij}^{(2)} = 0 \end{array}$$

$$\text{and } \frac{\partial J(\theta)}{\partial \theta_{ij}^{(1)}} = D_{ij}^{(1)} \quad \begin{array}{l} \text{can substitute this} \\ \text{in training algorithm. if} \\ \text{we do not want to} \\ \text{compute } \Delta_{ij}^{(1)} \end{array}$$

$$\delta_j^{(1)} = \frac{\partial J(\theta)}{\partial z_j^{(1)}} \text{ cont}(i)$$

$$\text{cost}(i) = y^{(i)} \log \theta_0(x^{(i)}) + (1-y^{(i)}) \log \theta_0(x^{(i)})$$

when  $\delta$  term changes, it changes the cost.

function also, function increases in efficiency.  
 they change weights.

all the values to be zero initially.

## Training a NN

- i) Pick a network architecture.
- ii) Randomly initialize weights
- iii) Implement forward propagation
- iv) Find cost function  $J(\theta)$
- v) Implement backprop to  $\frac{\partial J(\theta)}{\partial \theta^{(l)}}$
- vi) Use gd or any optimisation method with backprop to minimize  $J(\theta)$

Evaluating hypothesis:

training set : 60%

cross validation : 20%

test set : 20%

Select different hypothesis & train & find cross validation J error.

Select model with least cv error.

## Bias vs Variance

if degree of polynomial is low, then its bias problem.

if degree of polynomial is large, then its variance problem.

## Bias (Underfit)

$J_{\text{train}}$  is high.

$J_{\text{cv}} \approx J_{\text{train}}$

$J_{\text{cv}}$  is high, which is most likely due to underfitting.

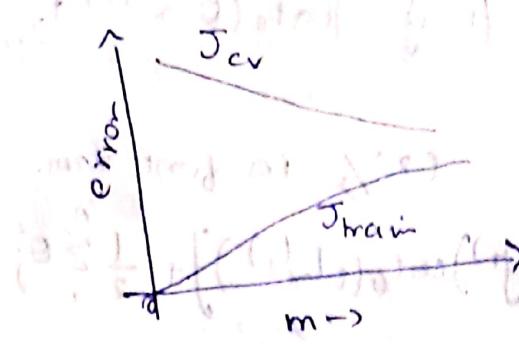
## Variance (Overfit)

$J_{\text{train}} \rightarrow \text{low}$

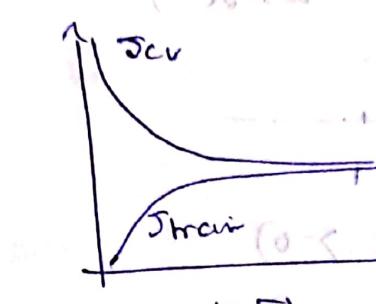
$J_{\text{cv}}$  is high

$J_{\text{cv}}$  is high, which is due to overfitting.

## Underfitting

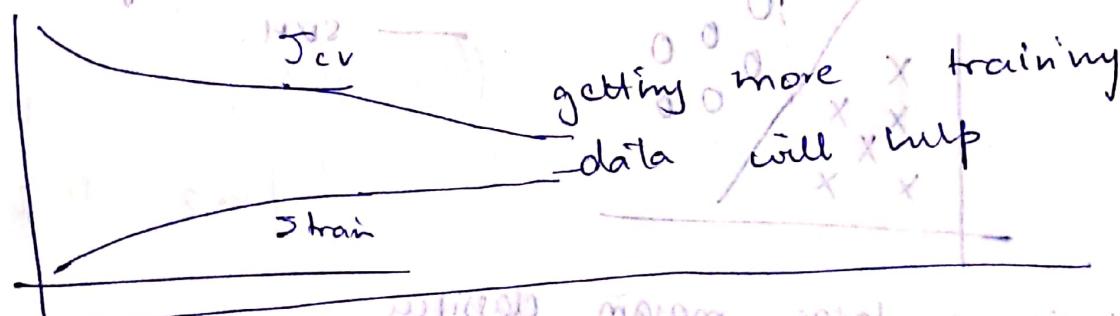


high bias



for high bias, getting more training data won't help.

high variance

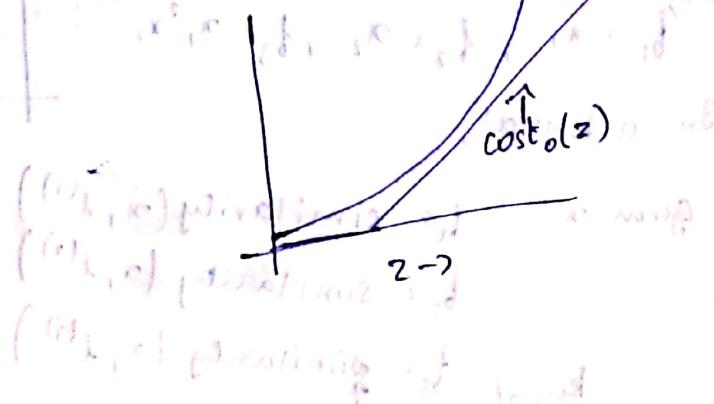
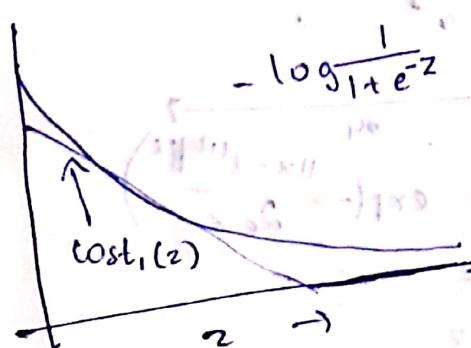


Support Vector Machine

when  $y \geq 1$ , we want  $\theta^T x \geq 0$

$y \geq 0$ , we want  $\theta^T x \geq 0$

$$-\log\left(\frac{1}{1+e^{-z}}\right)$$

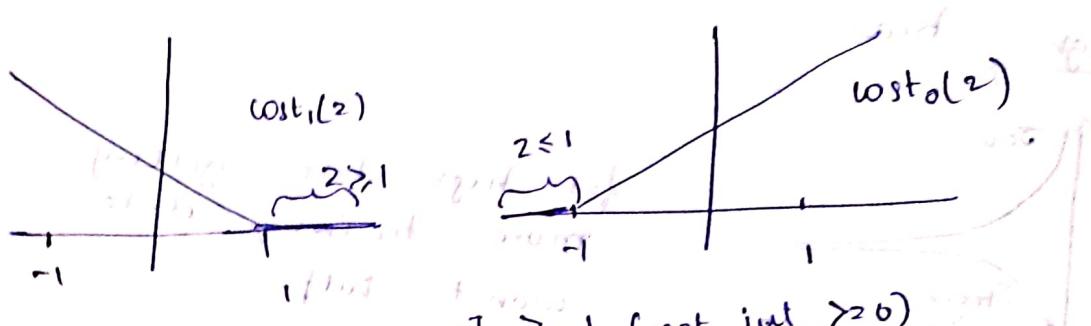


for support vector machine,

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{\lambda}{2m} \sum_{j=0}^n \theta_j^2$$

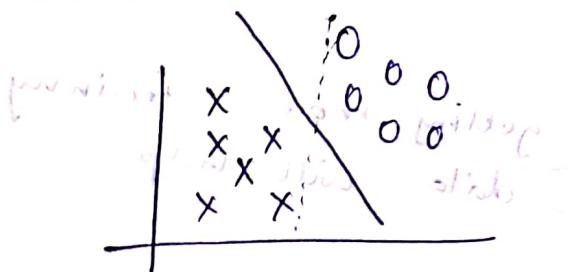
ignore  $m$ , introduce  $(2/2)$  in front term.

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=0}^n \theta_j^2$$



if  $y \geq 1$ , we want  $\theta^T x \geq 1$  (not just  $> 0$ )  
 if  $y \leq 0$ , we want  $\theta^T x \leq -1$  (not just  $< 0$ )

SVM decision boundary:



---> linear regression  
 —— SVM

it is a large margin classifier

### Kernels

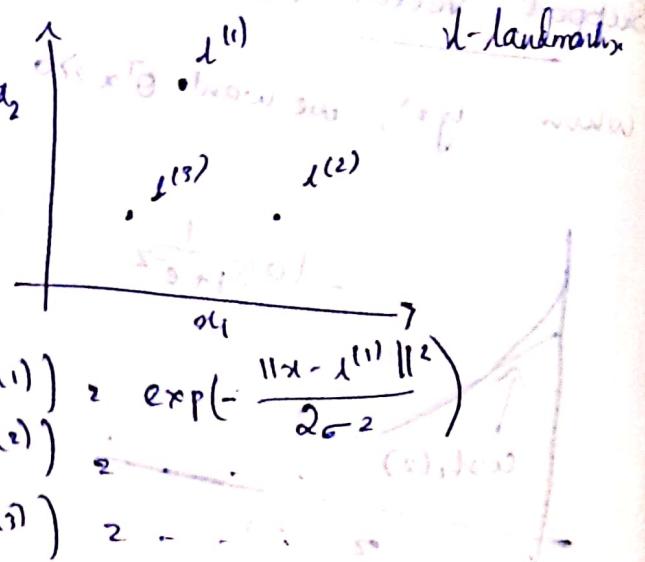
Given  $\theta = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^3 + \dots$   
 Here,  $f_1 = x_1, f_2 = x_2, f_3 = x_1^2 x_2, \dots$

In a kernel,

Given  $x$ :  $f_{1,2}$  similarity( $x, \lambda^{(1)}$ )  $\approx \exp(-\frac{\|x - \lambda^{(1)}\|^2}{2\sigma^2})$

$f_{2,3}$  similarity( $x, \lambda^{(2)}$ )  $\approx \dots$

$f_{3,4}$  similarity( $x, \lambda^{(3)}$ )  $\approx \dots$   
 Kernel function (gaussian kernel)  
 $\|\omega\| = \text{length of vector } \omega$



$$\text{if } \alpha \approx 1 \Rightarrow f_1(x) \approx 1 - \frac{\sigma^2}{2\alpha^2} \approx 1$$

If  $\alpha$  far from 1, then  $f_1(x)$  is small.

$$f_1(x) \exp\left(-\frac{\sigma^2}{2\alpha^2}\right) \approx 0$$

compute values of  $f_1(x)$  for different values of  $\alpha_1, \alpha_2, \alpha_3, \dots$

$\therefore$  In the plot we will have regions where  $y_1 > 0$  and regions where  $y_1 < 0$ . Hence we will get the decision boundary.

### SVM with kernels

given  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$ , choose  $\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(m)}$  satisfying  $\lambda_i \geq 0$

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

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

$$f_i = \begin{cases} 0 \\ \alpha_i \\ 1 \end{cases}$$

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

Predict  $y^{(i)}$  if  $\theta^T f_i > 0$  for all  $i$  and  $\theta^T f_i \leq 0$  for all  $i$ .

using SVM algorithm, replace  $\alpha_i$  by  $\alpha_i^{(i)}$  and  $\lambda_i$  by  $\lambda_i^{(i)}$ .

SVM parameter:  $C = \frac{1}{2} \sum \lambda_i$  or  $\lambda_i$  of max size in training

large  $C$ : low bias, high variance

low  $C$ : high bias, low variance

$\sigma^2$ : feature  $f$  vary more smoothly, High bias, low variance

Small  $\sigma^2$ : vary less smoothly, lower bias, higher variance.

To use SVM software package,

- specify (
- choose kernel

For SVM packages' optimisations to work properly, it should satisfy Mercer's theorem.

Other kernels:

polynomial kernel:  $k(x, l) = (x^T l)^2$  or  $(x^T l + 2)^4$

For Multiclass classification, one versus all method. Many kernels support this, otherwise follow one versus all method.

## UNSUPERVISED LEARNING

Clustering: data is given without it being classified. Algorithm groups it on its own.

### K-Means algorithm

i) cluster assignment:

In case you need to group it into 2<sup>nd</sup> group, initial two cluster centroids (random). It will go through each data and it will assign to one of the cluster centroid depending on which is nearer to it.



2) Now take mean of two groups & move the centroid to mean value.

Repeat step ① & ② for the new centroids, till the centroids won't change.

$k \rightarrow$  no. of clusters.

Randomly initialise  $K$  cluster centroids  $\mu_1, \mu_2, \mu_3, \dots, \mu_k$

Repeat {  
for  $i = 1$  to  $m$   
(i) = index [1 to  $K$ ] of cluster centroid closest  
 $(\|\mathbf{x}^{(i)} - \mu_k\|)$  to  $\mathbf{x}^{(i)}$  assigned to k}

for  $k = 1$  to  $K$  (loop)  
assign  $\mu_k = \text{avg}(\text{mean})$  of points assigned to  $k$ .

if max } least or biggest in std values of  
{  $\mu_k$  } then stop else go to step 10

optimisation objective:

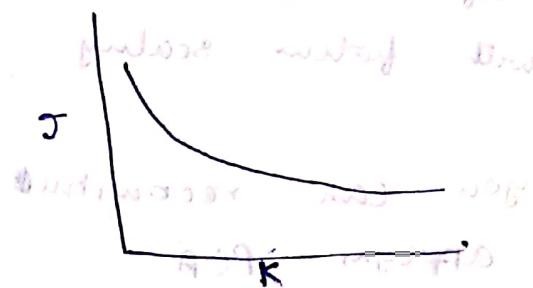
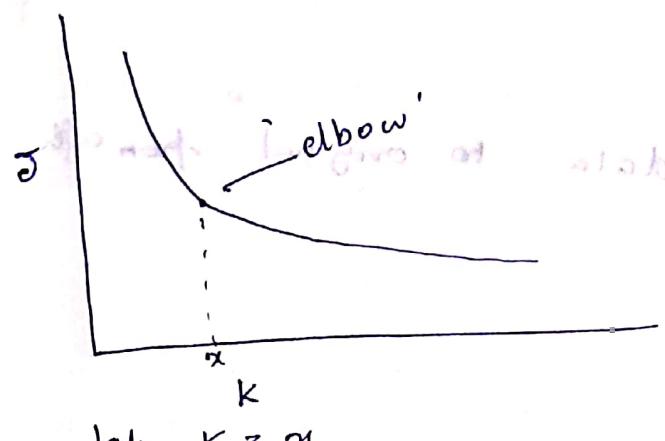
$$J(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i - \mu_i\|^2$$

minimise  $J(\cdot)$

In step 1, need to minimise  $J$  for  $c$  values keeping  
 $\mu$  constant. In step 2, need to minimise  $J$   
for  $\mu$  values.

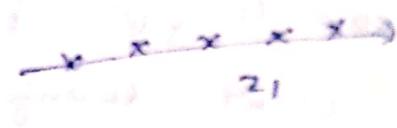
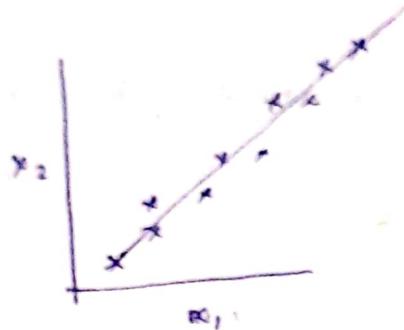
Choosing value of  $K$

elbow method



difficult to choose  $K$

Data compression  $20 \rightarrow 10$



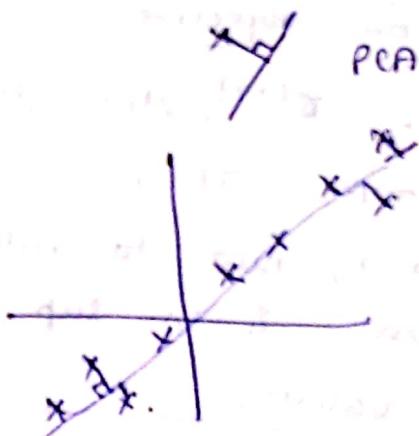
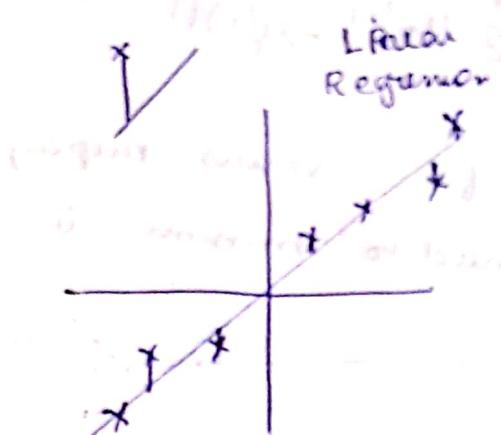
Same can be done for 3D to 2D.

When large no. of features are given, you can compute combination of them & reduce no. of features for better visualisation.

(Ch 10)

### PCA - Principal Component Analysis

It tries to find a lower dimensional surface onto which data is projected so that sum of square from data to projected plan / line is minimum.

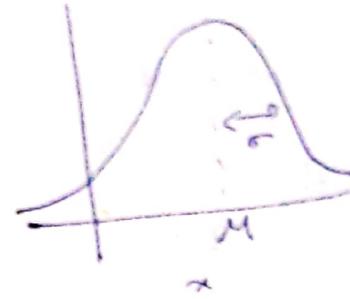


Perform mean normalisation before PCA along with feature scaling

You can reconstruct data to original from after applying PCA

## Anomaly Distribution

is distributed gaussian with mean  $\mu$ , variance  $\sigma^2$



$$p(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

### Density estimation

$$p(x) = p(x_1, \mu_1, \sigma_1^2) p(x_2, \mu_2, \sigma_2^2) \dots p(x_n, \mu_n, \sigma_n^2)$$

$$\Rightarrow \prod p(x_j, \mu_j, \sigma_j^2)$$

### Anomaly detection

choose  $x_j$  that might be indicative of anomalous examples.

Fit  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum x_j^{(1)}$$

$$\sigma_j^2 = \frac{1}{m} \mathbb{E} \sum (x_j^{(1)} - \mu_j)^2$$

for a new  $x$ ,

$$p(x) = \prod p(x_j, \mu_j, \sigma_j^2)$$

Anomaly if  $p(x) < \epsilon$

### Algorithm evaluation

Fit  $p(x)$  on training set

on CV / test set, predict

$$y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases}$$

Possible evaluation metrics based on confusion matrix

- true positive, false positive, true negative, false negative
- accuracy, recall, F1-score

Anomaly detection vs Supervised learning

- Small no of positive anomaly examples.
- Fraud detection
- Manufacturing defect
- Email spam classification
- Weather prediction

Choosing features for anomaly detection

- better when selected feature have proper value
- if a feature won't give out proper value, transform it. Eg  $x_i \rightarrow \log(x_i + 1)$

Predicting movie problem.

train feature with user input

collaborative filtering algorithm

Dealing with large datasets. Stochastic GD

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

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

Randomly shuffle dataset

Here we are looking at training example  
in modifying parameters

while in Batch GD (all at once) whole  
feature sum is made in parameters changed.

Mini-Batch GD

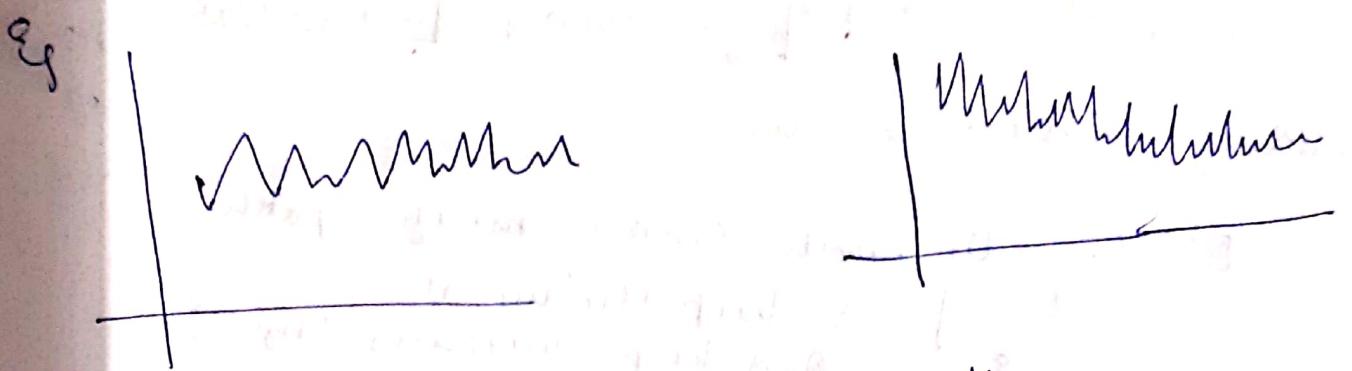
In stochastic, we use only 1 example in each iteration.

In mini-Batch, use b examples in each iteration

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

## Checking for convergence.

Cost fn vs No of iterations - curve isn't smooth.



Can slowly decrease & over time

? const

iteration no + const<sup>2</sup>

## Map Reduce

Batch GD:  $\theta_j = \theta_j - \frac{\alpha}{400} \sum_{i=1}^{400} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$

split training set into several sub sets.

Machine 1: first 100

$$\text{temp}_j^{(1)} = \sum_{i=1}^{100} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Machine 2: next 100

$$\text{temp}_j^{(2)} = \sum_{i=101}^{200} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Combine

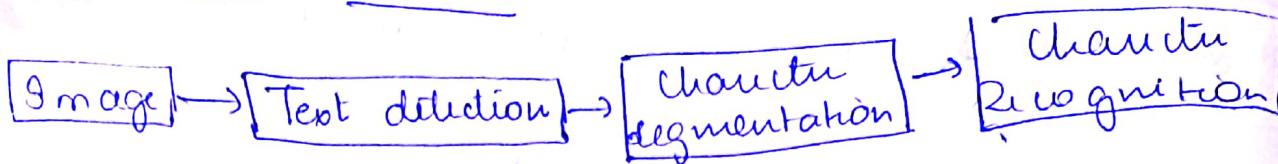
$$\theta_j = \theta_j - \frac{\alpha}{400} (\text{temp}_j^{(1)} + \text{temp}_j^{(2)} + \text{temp}_j^{(3)}, \text{temp}_j^{(4)})$$

In single computers with multiple core machines,  
it can split & computed.

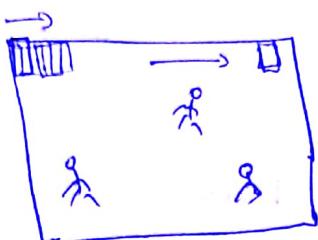
Photo

OCR

Pipeline



Sliding window detection.



take small image patch  
 & keep sliding it.  
 And keep increasing size of  
 patch.

Decision - Tree.

Data is continuously split according to a certain parameter.

two entities;

decision nodes → data is split leaves → final outcomes.

two types:

classification trees

decision variable is categorical

Regression trees

outcome variable is continuous

algorithm for decision tree

ID3 algorithm, Iterative Dichotomiser 3

Basic definitions

Entropy: for a finite S, it is the measure of uncertainty or randomness in data

Information gain: (Kullback-Leibler divergence)  $I(S, A)$  for a set  $S$  is the effective change in entropy after deciding on a particular attribute  $A$ . It measures a relative change in entropy with respect to the independent variables.

$$I(S, A) = H(S) - H(S, A)$$

$$\Rightarrow H(S) = \sum P(x) * H(x)$$

$I(S, A)$  is the information gained by applying feature  $A$ .

Steps:

- Create root node for this tree.
- If all examples are positive, return leaf node.
- If all examples are negative, return leaf node.
- Calculate the entropy of the current state  $H(S)$ .
- For each attribute, calculate the entropy with respect to the attribute 'x' denoted by  $H(S, x)$ .
- Select the attribute which has maximum value of  $I(S, x)$ .
- Remove the attribute that offers highest  $I(S, x)$  from set of attributes.
- Repeat until we run out of all attributes or the decision tree has all leaf nodes.
- Select the feature as root node which has max information gain.
- And so on for following nodes.