

Reference Books :

- 1.) Pattern Recognition & Machine learning by Christopher M. Bishop
- 2.) Deep learning by Ian Good Fellow, Yoshua Bengio , and Aaron Courville
- 3.) Machine Learning By Andrew Ng (standford)
- 4.) Machine learning by Stanley chan (Purdue)

Mid Term - 30

Quizzes (2) - 20

End Term - 50

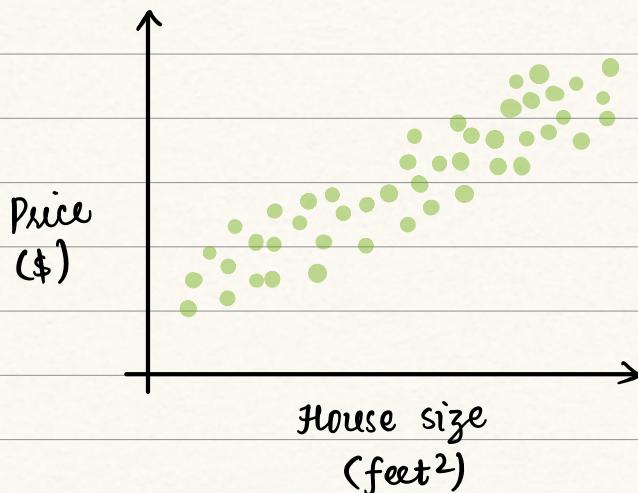
Theory class + Demonstration of Machine Learning (ML)
Algorithms on Python Platform

→ Mathematically challenging - Knowledge of statistics, Basic Probability & Linear Algebra.

* Supervised Machine Learning —

→ Labelled data set which is having both feature as well values of corresponding feature .

Example : House size and Price in \$.



Q.) What is the goal based on this dataset ?

→ Predicting the house price close to the actual price.

Q) Given the dataset like we seen before, how can we learn to predict the prices of the other houses which might or might not be present in the dataset?

$x^{(i)}$ → denotes the 'input' variable in the dataset.

This superscript (i) denotes the i^{th} training example.

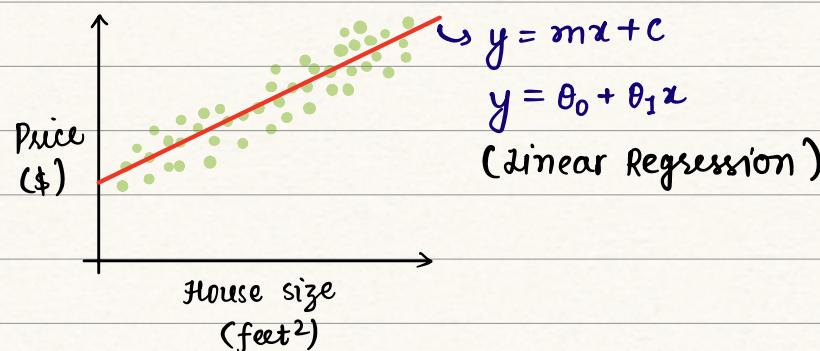
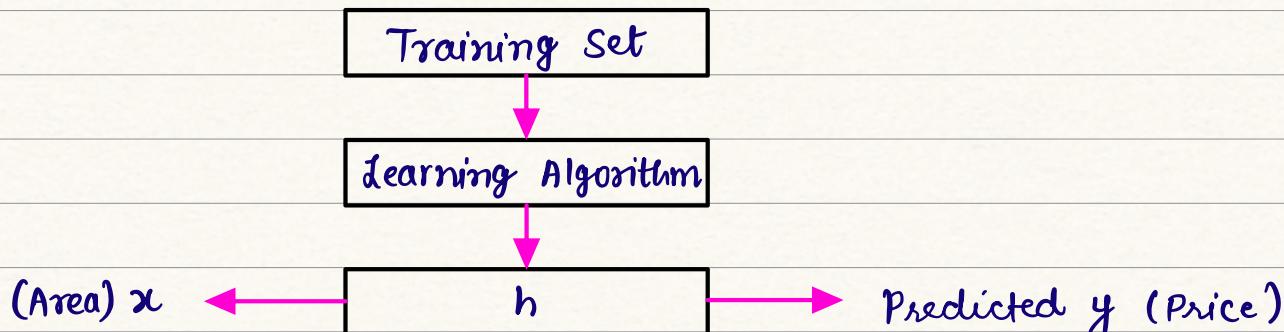
$y^{(i)}$ → denotes the 'output' variable that we are trying to predict.

$(x^{(i)}, y^{(i)})$ → This pair is called i^{th} training example in the dataset.

$\{(x^{(i)}, y^{(i)}) ; i = 1, 2, \dots, N\}$

This is the dataset.

Given a training set, we need to learn a function $h: x \rightarrow y$ so that $h(x)$ is a 'good' predictor for corresponding values of y . { $h \rightarrow$ hypothesis function}



$$h(x) = \theta_0 + \theta_1(x) = [\theta_0 \ \theta_1] \begin{bmatrix} 1 \\ x \end{bmatrix} = \theta^T x$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \quad \theta^T = [\theta_0 \ \theta_1]$$

Q.) Given a training dataset, choose the parameters θ_0 and θ_1 ?

We should choose the values of θ_0 and θ_1 such that

$\underbrace{h(x^{(i)})}_{\text{Predicted } y^{(i)}} = \theta_0 + \theta_1 x^{(i)}$ is close to each $y^{(i)}$ (True value).

$$\text{cost function } J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2$$

Mean squared Error (MSE)

Our aim to find a way to compute the optimum values of θ_0 and θ_1 that minimize MSE, $J(\theta)$.

Gradient Descent Algorithm

$$F = x_1 f_1 + x_2 f_2 + \dots + x_n f_n$$

$$\nabla F = \begin{bmatrix} \frac{\partial F}{\partial f_1} \\ \frac{\partial F}{\partial f_2} \\ \vdots \\ \frac{\partial F}{\partial f_n} \end{bmatrix}$$

* starts with some initial guess of θ_0 & θ_1 .
 * Algorithm repeatedly changes the values of θ_0 and θ_1 till $J(\theta)$ converges to the minimum value.

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

Assignment operator
 Learning rate

This update is simultaneously done for θ_0 & θ_1 .

* learning rate should be optimum, it should not be too high and too low.

$$n=1;$$

$$J(\theta) = \frac{1}{2} (h(x) - y)^2$$

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = (\theta_0 + \theta_1 x - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 + \theta_1 x - y)$$

$$= (\theta_0 + \theta_1 x - y) \cdot x_j$$

$$\left\{ x_0^{(i)} = 1, x_j^{(i)} = x^{(i)} \right\}$$

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = (h(x) - y) \cdot x_j$$

$$\theta_j := \theta_j + \alpha (y - h(x)) \cdot x_j \quad \text{for one training}$$

example

$$\theta_j := \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) \cdot x_j^{(i)}$$

Repeat until convergence of $J(\theta)$.

$$\theta_j := \theta_j + \frac{\alpha}{n} \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) \cdot x_j^{(i)}$$

for every $j = 0, 1$

$$\theta_0 := \theta_0 + \frac{\alpha}{n} \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) \cdot x_0^{(i)}$$

$$\theta_1 := \theta_1 + \frac{\alpha}{n} \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) \cdot x_1^{(i)}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \quad x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \end{bmatrix}$$

Subject:

Batch
Gradient Descent

Convex function

30 /07/2025

$$\theta_j := \theta_j + \frac{\alpha}{n} \sum_{i=1}^n \underbrace{(y^{(i)} - h(x^{(i)})) \cdot x^{(i)}}_{\text{Error}} \quad \{j = 0, 1, 2\}$$

Q: what will happen if the error is very small & very high?

Weights θ_0 & θ_1 → optimized (in case of very small)

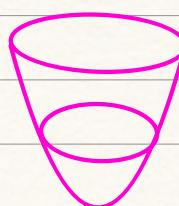
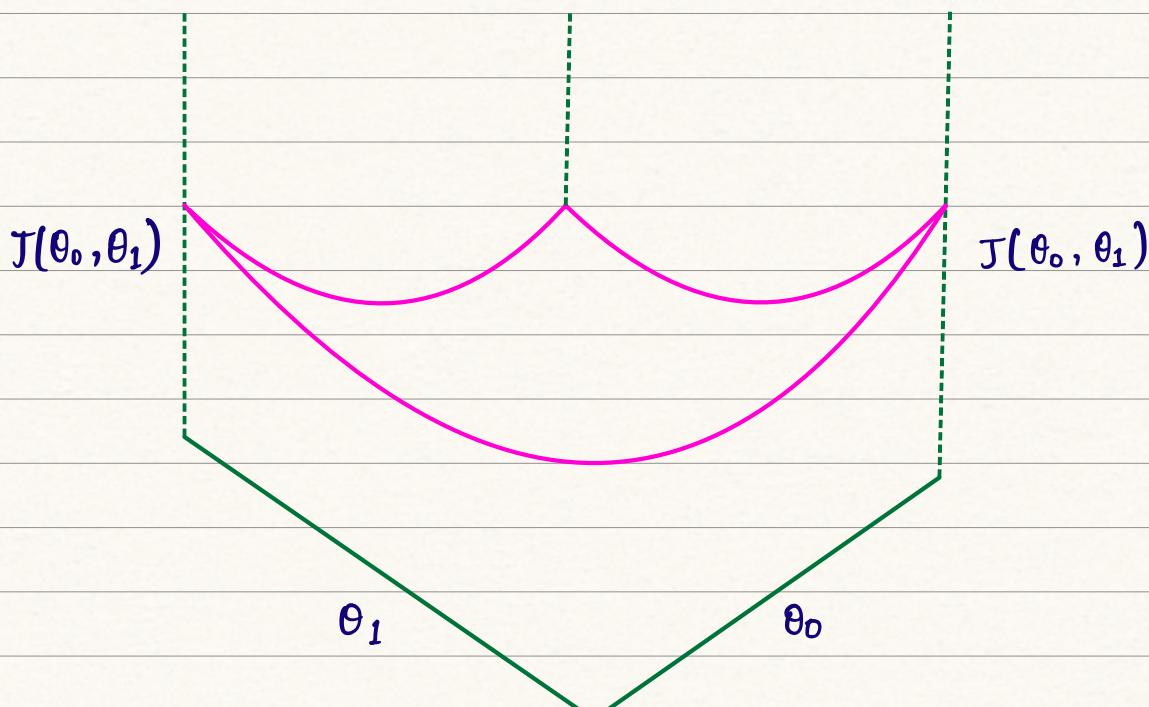
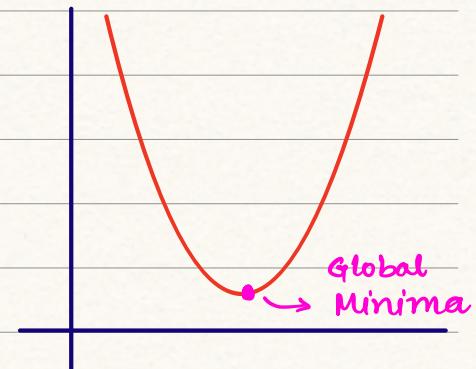
Weights θ_0 & θ_1 → not optimized (in case of very high)

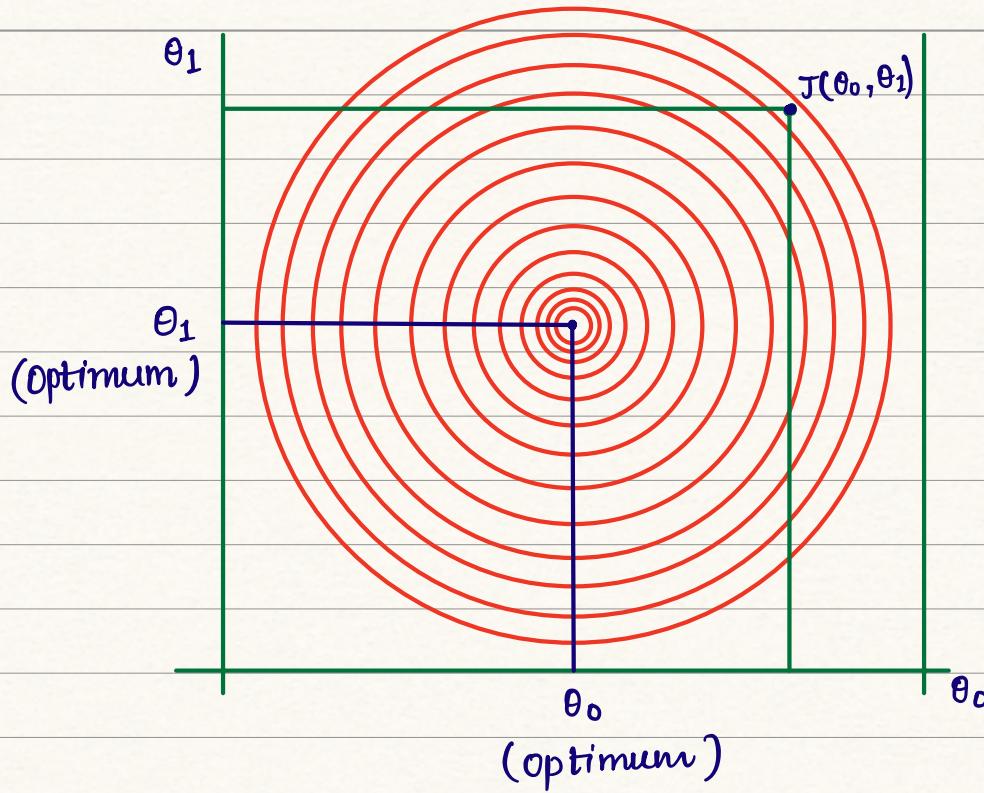
when we have a very large dataset, we don't use the batch gradient descent.

* Convex function : —

with GDA, $J(\theta)$ achieves the global minimum value.

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





* Stochastic Gradient Descent (Incremental Gradient Descent) :-

$$\theta_j := \theta_j + \alpha (y^{(i)} - h(x^{(i)})) \cdot x_j^{(i)} ; j = 0, 1 \\ \text{for } i = 1 \text{ to } n.$$

In the stochastic gradient descent algorithm, whenever we encounter the Training example, we update the parameters (θ_0, θ_1) according to the gradient of the error with respect to that single training example only.

* With the stochastic gradient Descent Algorithm, we might never achieve the actual minimum value of the cost function and the achieved cost function with this algorithm might be oscillating around the global minimum value.

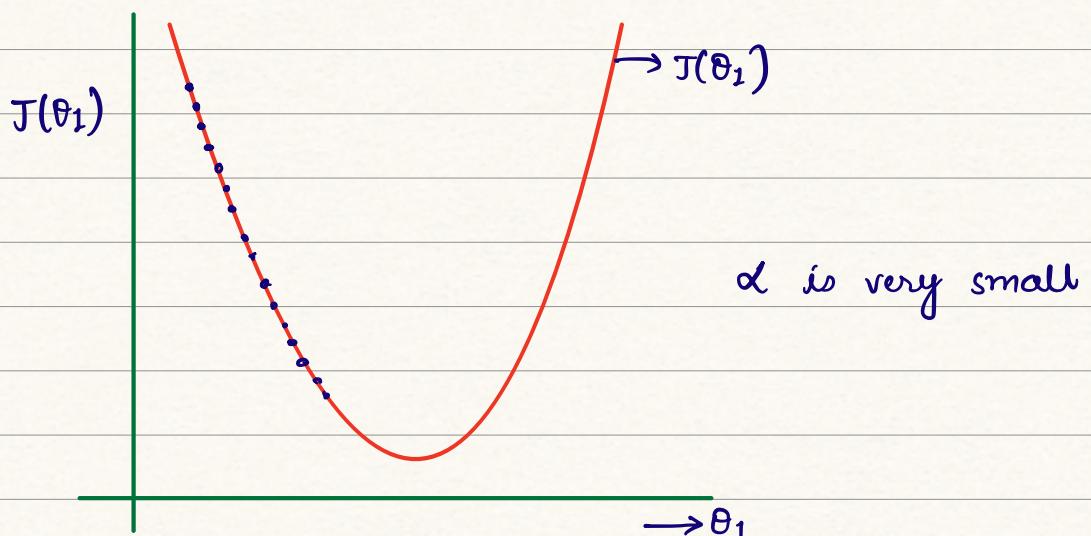
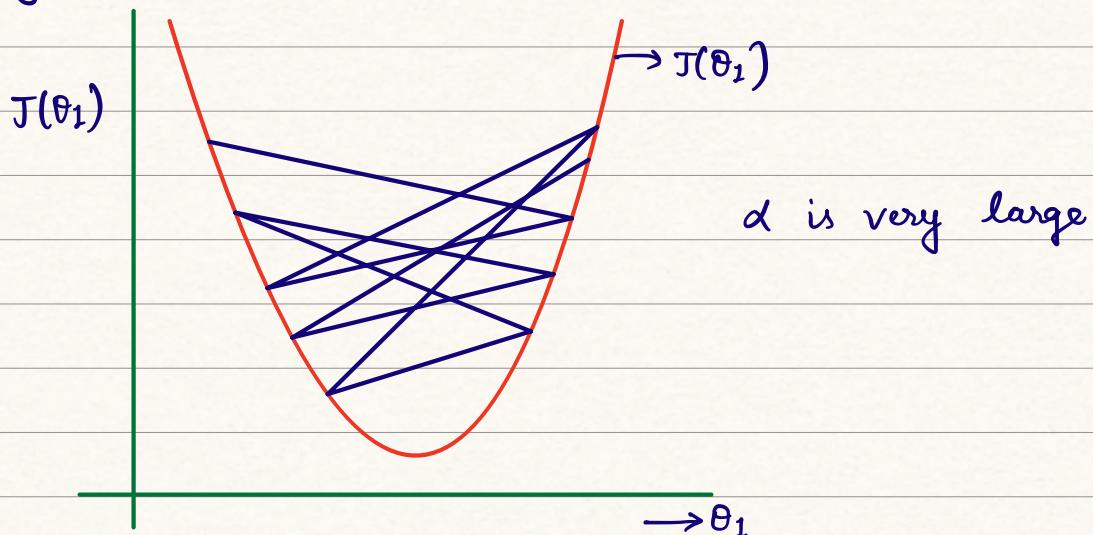
* Still, although it fails to achieve the global minimum value, it is found that the cost function value obtained is approximately very close to the global minimum value & hence we can still use that approximated achieved cost function value.

Subject:

/ /

- * When the size of the training set is very large, stochastic gradient descent is preferred over the Batch Gradient Descent.
- * The parameter α (learning rate) should be chosen very carefully for the better convergence of the algorithm.
- * If the α is chosen very large, Gradient Descent might fail to converge.

$$y = \theta_1 x$$



- * If α is very small, GDA would take very much time to achieve the minimum value of $J(\theta_1)$.

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

Multivariate

 $x_1 \rightarrow$ size of house (ft^2)

Linear

 $x_2 \rightarrow$ no. of bedrooms

Regression

 $x_3 \rightarrow$ no. of gym $x_4 \rightarrow$ garage or not

$$\theta_j := \theta_j + \frac{\alpha}{n} \sum_{i=1}^n (y^{(i)} - h(x^{(i)})). x_j^{(i)}$$

$j = 0, 1, 2, 3, 4$

* In general, we can have ' d ' number of features.

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d$$

$$\begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_d \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

★ Normal Equations :—

* GDA gives one way to find the parameters/weights (θ_0, θ_1) in order to minimize the mean squared error (MSE).

Let's study another way to find the optimum values of (θ_0, θ_1) so as to minimize $J(\theta)$.

Multi-features data set regarding the housing price features (living area (ft^2), No. of bedrooms, Price (\$))

x_0	Living area (feet 2)	# of bedrooms	Price (\$10,000)	
1	2104	3	400	Training Example 1
1	1600	3	330	
1	2400	3	369	

Subject:

$$\theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)} \rightarrow \text{First training example}$$

d input features

no. of training examples $\begin{cases} 1 \\ 2 \\ \vdots \\ n \end{cases}$

$$x_0 \quad x_1 \quad x_2 \quad \dots \quad x_d$$

$$\text{size} = n \times (d+1)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_d \end{bmatrix} (d+1) \times 1$$
$$X = \begin{bmatrix} - x^{(1)\top} - \\ - x^{(2)\top} - \\ - x^{(3)\top} - \\ - x^{(4)\top} - \\ - \vdots - \\ - x^{(n)\top} - \end{bmatrix} n \times (d+1)$$

$$x\theta - \vec{y}$$

$$\begin{bmatrix} - x^{(1)\top} - \\ - x^{(2)\top} - \\ - x^{(3)\top} - \\ - x^{(4)\top} - \\ - \vdots - \\ - x^{(n)\top} - \end{bmatrix} n \times (d+1) \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_d \end{bmatrix} (d+1) \times 1 \quad - \quad \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ \vdots \\ y_5 \end{bmatrix} n \times 1$$

$$\underbrace{\theta_0 + \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)} + \dots + \theta_d x_d^{(1)}}_{h(x^{(1)})} \quad ? \quad \text{first training example}$$

$$\Rightarrow h(x^{(i)}) - y^{(i)}$$

Subject:

$$\begin{bmatrix} h(x^{(1)}) \\ h(x^{(2)}) \\ \vdots \\ h(x^{(n)}) \end{bmatrix}_{nx1} - \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}_{nx1} = \begin{bmatrix} h(x^{(1)}) - y^1 \\ h(x^{(2)}) - y^2 \\ \vdots \\ h(x^{(n)}) - y^n \end{bmatrix}_{nx1}$$

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2 \rightarrow \text{MSE}$$
$$= \frac{1}{2n} \Delta^T \Delta$$

where $\Delta = \begin{bmatrix} h(x^{(1)}) - y^1 \\ h(x^{(2)}) - y^2 \\ \vdots \\ h(x^{(n)}) - y^n \end{bmatrix}$

$$= \frac{1}{2n} (x\theta - \bar{y})^T (x\theta - \bar{y})$$

Given $c \in \mathbb{R}^n$, the differential function $f(x) = c^T x ; x \in \mathbb{R}^n$.

$$f(x) = c^T x$$
$$= c_1 x_1 + c_2 x_2 + \dots + c_n x_n$$

$$c = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix}_{nx1}$$

$$c^T x = [c_1 \ c_2 \ \dots \ c_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$c^T x = c_1 x_1 + c_2 x_2 + \dots + c_n x_n$$

Gradient ($\nabla_x f(x)$) = $\begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$

Subject:

/ /

For the sake of understanding, take $d=2$;

$x_0 \quad x_1 \quad x_2 \rightarrow$ Suppose two training examples

$$x = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \end{bmatrix}_{2 \times 3}$$

$$x^T = \begin{bmatrix} 1 & 1 \\ x_1^{(1)} & x_1^{(2)} \\ x_2^{(1)} & x_2^{(2)} \end{bmatrix}_{3 \times 2}$$

$$x^T x = \begin{bmatrix} 1 & 1 \\ x_1^{(1)} & x_1^{(2)} \\ x_2^{(1)} & x_2^{(2)} \end{bmatrix}_{3 \times 2} \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \end{bmatrix}_{2 \times 3}$$

$$x^T x = \begin{bmatrix} 2 & x_1^{(1)} + x_1^{(2)} & x_2^{(1)} + x_2^{(2)} \\ x_1^{(1)} + x_1^{(2)} & (x_1^{(1)} + x_1^{(2)})^2 & x_1^{(1)} x_2^{(1)} + x_1^{(2)} x_2^{(2)} \\ x_2^{(1)} + x_2^{(2)} & x_1^{(1)} x_2^{(1)} + x_1^{(2)} x_2^{(2)} & (x_2^{(1)} + x_2^{(2)})^2 \end{bmatrix}$$

$$x^T x = B = B^T$$

For n training examples and having ' d ' no. of input features (excluding $x_0 = 1$). The matrix $x^T x$ is always symmetric.

Suppose that $A = A^T$ (symmetric matrix).

$$\nabla_x (x^T A x) = \boxed{\quad}$$

$$A = \begin{bmatrix} a_1 & a_2 \\ a_2 & a_1 \end{bmatrix}$$

$$x^T A x = [x_1 \quad x_2] \begin{bmatrix} a_1 & a_2 \\ a_2 & a_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$x^T A x = a_1 x_1^2 + a_2 x_1 x_2 + a_2 x_1 x_2 + a_1 x_2^2$$

$$\nabla_x(x^T A x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2a_1 x_1 + a_2 x_2 + a_2 x_2 \\ a_2 x_1 + a_2 x_1 + 2a_1 x_2 \end{bmatrix}$$

$$= 2 \begin{bmatrix} a_1 & a_2 \\ a_2 & a_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\nabla_x(x^T A x) = 2Ax$$

$$\begin{aligned} J(\theta) &= \frac{1}{2n} (x\theta - y)^T (x\theta - y) \\ &= \frac{1}{2n} ((x\theta)^T - y^T)(x\theta - y) \\ &= \frac{1}{2n} [(x\theta)^T x\theta - (x\theta)^T y - y^T x\theta + y^T y] \quad \xrightarrow{\frac{\partial J(\theta)}{\partial \theta} = 0} \\ &= \frac{1}{2n} [\theta^T x^T x\theta - y^T x\theta - y^T x\theta + y^T y] \\ ((x\theta)^T y)^T &= y^T ((x\theta)^T)^T = y^T x\theta \text{ (scalar)} \\ &= \frac{1}{2n} [\theta^T x^T x\theta - 2y^T x\theta] \quad []_{mxn} \rightarrow \text{matrix} \\ &= \frac{1}{2n} [\theta^T x^T x\theta - 2(x^T y)^T \theta] \quad []_{mx1} \rightarrow \text{vector} \end{aligned}$$

$$\nabla_\theta J(\theta) = \frac{1}{2n} [2x^T x\theta - 2x^T y]$$

$$\nabla_\theta J(\theta) = 0$$

$$\begin{aligned} \Rightarrow x^T x\theta &= x^T y \\ \Rightarrow \theta &= (x^T x)^{-1} x^T y \end{aligned}$$

$$\left\{ \begin{array}{l} \nabla_x x^T A x = 2Ax \\ b \in \mathbb{R}^n \quad x \in \mathbb{R}^n \\ f(x) = b^T x \\ \nabla_x f(x) = b \\ \nabla_\theta f(\theta) = (x^T y)^T \theta \\ \quad \quad \quad = x^T y \end{array} \right.$$

Subject:

$$\theta = \underbrace{(x^T x)}_{\text{Pseudo inverse of } X}^{-1} x^T y$$

$$\left\{ \begin{array}{l} (x^T x)^{-1} x^T \\ = x^{-1} (x^T)^{-1} x^T \\ = x^{-1} \end{array} \right.$$

Q: How dependent input features can be a problem?

$$(x^T x)^{-1} x^T y$$

↳ can create a problem when
the features in the data set are linearly dependent.



Determinant = 0

$$A^{-1} = \frac{\text{Adjoint of } A}{|A|}$$

$|A| \rightarrow$ here it will create a problem

Gradient Descent Algorithm	Normal Equation
* Need to choose the learning rate α .	* No need to choose α .
* Need many iterations.	* Don't need to do iteration.
* Works well even the data set is very large.	* Need to compute $(x^T x)^{-1}$ $O(d^3) \leftarrow$ * Works slow when d is large.

* In general, GDA is the recommended method for finding the parameters $(\theta_0, \theta_1, \theta_2, \dots, \theta_d)$.

★ Gaussian Distribution —

The PDF of a Gaussian Distributive random variable x having the parameters μ (mean) and σ^2 (variance) is

Subject:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} ; -\infty < x < \infty$$

$$x \sim N(\mu, \sigma^2)$$

read as x is modeled as Normal or Gaussian distribution function.

Standard Normal —

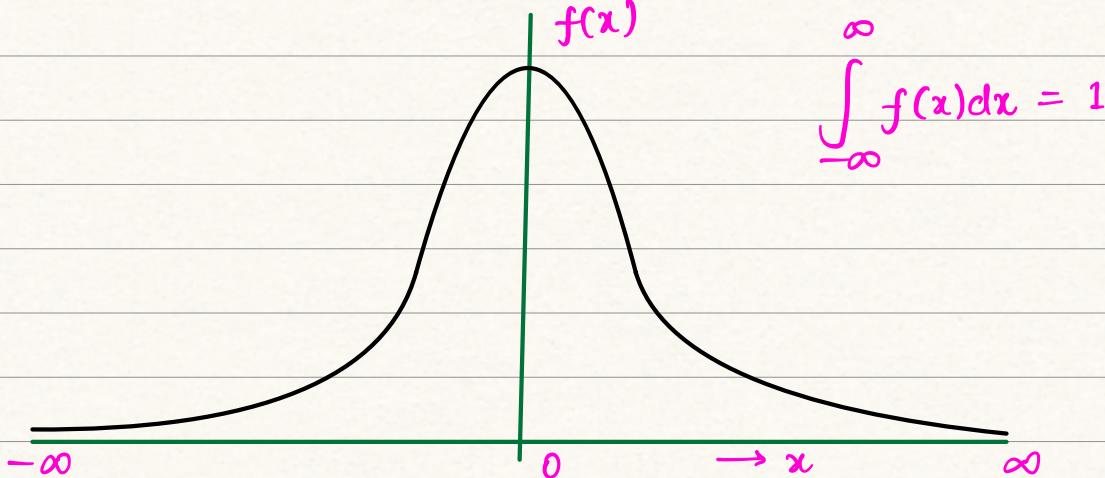
$$\mu = 0 \quad \text{and} \quad \sigma^2 = 1$$

$$x \sim N(0, 1)$$

→ expectation operator

$$\mu = E[x]$$

$$\sigma^2 = E[x^2] - (E[x])^2$$



$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$P(1 < x \leq 2) = \int_1^2 f(x) dx$$

* Suppose that we have five random values from a normal distribution, but we don't know which parameters of normal distribution it belongs to i.e., we need to find the parameters μ and σ^2 .

$$y = \underbrace{\{48.7, 50.3, 47, 51, 48\}}$$

It's not coming from standard normal. $N(0, 1)$.

Subject:

Q.) what is the likelihood of y generated from $N(0, 1)$?
Almost zero.

Q.) What is the likelihood of y generated from $N(50, 1)$?
Much much higher

* This is the sample example for the intuition of estimation of the parameters μ & σ^2 . This approach in statistics is called maximum likelihood estimation (MLE).

↳ Find the parameters that maximize the likelihood of generating the data we observe.

$$x \sim N(\mu, \sigma^2)$$

$$f_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right); -\infty < x < \infty$$

In this set up, we have a complete information about the parameters μ & σ^2 . But what if the opposite is true - we know the outcome of the random variable x , but we do not know the parameters that generated it.

↳ likelihood function

$$L = (\underbrace{\mu, \sigma^2}_{\text{Function of unknown parameters}} | x) \xrightarrow{\text{data}}$$

Function of unknown parameters

$$L = \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)}_{\text{But we don't know } \mu \text{ & } \sigma^2 \text{ here.}}$$

This pdf because we know x is coming from Gaussian Distribution.

Subject:

x_1, x_2, \dots, x_n are independent & identically distributed (IID) - Random variables.

$$f(x_1, x_2, \dots, x_n) = f(x_1) \underbrace{f(x_2) \dots f(x_n)}_{\text{Joint pdf}} = \prod_{i=1}^n f(x_i)$$

All are identical

$$= (f(x))^n$$

$x_1, x_2, \dots, x_n \rightarrow$ IID samples from $N(\mu, \sigma^2)$

$$\begin{aligned} L(\mu, \sigma^2 | x) &= \prod_{i=1}^n f(x_i | \mu, \sigma^2) \\ x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \end{aligned}$$

Now we will find the parameters μ & σ^2 which maximizes the likelihood function $L(\mu, \sigma^2 | x) \rightarrow$ hence name is MLE.

$$L(\mu, \sigma^2 | x) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Taking log on both the sides;

$\ln(L(\mu, \sigma^2 | x)) \rightarrow$ log likelihood function

$$\ln(L(\mu, \sigma^2 | x)) = \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$\ln\left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)\right) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

Subject:

/ /

$$\ln(L(\mu, \sigma^2 | x)) =$$

$$\frac{\partial}{\partial \mu} \ln(L(\mu, \sigma^2 | x)) = 0$$

$$\sum_{i=1}^n (x_i - \mu) = 0$$

$$\Rightarrow \hat{\mu}_{MLE} = \left\{ \frac{1}{n} \sum_{i=1}^n x_i \right\}$$

sample mean

$$\frac{\partial}{\partial \sigma^2} \ln(L(\mu, \sigma^2 | x)) = 0$$

$$\Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\{\mu = \hat{\mu}_{MLE}\}$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$$

Q) While doing the linear regression, why did we chose the mean squared error (MSE) as a cost function?

$$x_0^{(1)} \quad x_1^{(1)} \quad x_2^{(1)} \quad y^{(1)} \quad \theta^T x^{(1)}$$

$$x_0^{(2)} \quad x_1^{(2)} \quad x_2^{(2)} \quad y^{(2)} \quad \theta^T x^{(2)}$$

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

$$y^{(2)} - \theta^T x^{(2)} = e^{(2)}$$

Error

Assumptions :

* Each $e^{(i)}$; $i \in \{1, 2, 3, \dots, n\}$

$n = \#$ of training examples

$$e^{(i)} \sim N(0, \sigma^2)$$

Independent & Identically Distributed

$$f(e^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(e^{(i)})^2}{2\sigma^2}\right); -\infty < e^{(i)} < \infty$$

$i \in \{1, 2, \dots, n\}$

$$\Rightarrow y^{(i)} = \theta^T x^{(i)} + e^{(i)}$$

For a given $x^{(i)}$ and parameterized by θ written as
 $y^{(i)} | x^{(i)}; \theta$

$$y^{(i)} | x^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$$

$$E[x+y] = E[x] + E[y]$$

$x \rightarrow$ Deterministic

\hookrightarrow is not random

$$\begin{aligned} E[y^{(i)}] &= E[\theta^T x^{(i)} + e^{(i)}] \\ &= E[\theta^T x^{(i)}] + E[e^{(i)}] \\ &= \theta^T x^{(i)} + 0 \\ &= \theta^T x^{(i)} \end{aligned}$$

$$E[x] = x$$

$$\text{Var}(x) = E[x^2] - (E[x])^2$$

$$y^{(i)} = \theta^T x^{(i)} + e^{(i)}$$

$$\text{Var}(y^{(i)}) = E[(\theta^T x^{(i)} + e^{(i)})^2] - (E[\theta^T x^{(i)} + e^{(i)}])^2$$

$$\Rightarrow \text{Var}(y^{(i)}) = E[(\theta^T x^{(i)})^2 + (e^{(i)})^2 + 2\theta^T x^{(i)} e^{(i)}] - (E[\theta^T x^{(i)}] + E[e^{(i)}])^2$$

$$\begin{aligned} \Rightarrow \text{Var}(y^{(i)}) &= (\theta^T x^{(i)})^2 + \sigma^2 + 0 - (\theta^T x^{(i)} + 0)^2 \\ &= \cancel{(\theta^T x^{(i)})^2} + \sigma^2 - \cancel{(\theta^T x^{(i)})^2} \\ &= \sigma^2 \end{aligned}$$

$$y^{(i)} \sim N(\theta^T x^{(i)}, \sigma^2)$$

Subject:

$$f(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

Joint distribution of $y^{(i)}$: $f(y^{(1)}, y^{(2)}, \dots, y^{(n)})$

$$f(y^{(1)}, y^{(2)}, \dots, y^{(n)} | x; \theta) = f(y^{(1)} | x^{(1)}; \theta) f(y^{(2)} | x^{(2)}; \theta) \dots f(y^{(n)} | x^{(n)}; \theta)$$

$$\begin{aligned} f(y^{(1)}, y^{(2)}, \dots, y^{(n)} | x; \theta) &= \prod_{i=1}^n f(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= L(\theta | x, y) \quad \text{likelihood function} \end{aligned}$$

$$\ln L(\theta | x, y) = \ln \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \right)$$

$$\log(abc) = \log(a) + \log(b) + \log(c)$$

$$\ln L(\theta | x, y) = n \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

$$\text{Maximizing } \ln L(\theta | x, y) \rightarrow \text{minimizing } \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

* least square regression corresponds to the maximum likelihood estimate of θ .

★ Polynomial Curve Fitting :—

Training set having N observations ; $x = [x_1, x_2, \dots, x_N]^T$
and the corresponding outcomes to the observations are

$$\mathbf{t} = [t_1, t_2, \dots, t_N]^T$$

$N=10$ observations ;

$$\begin{array}{lcl} x_1 & = & 0.01 \\ x_2 & = & 0.005 \\ \vdots & = & \dots \\ x_{10} & = & 0.2 \end{array}$$

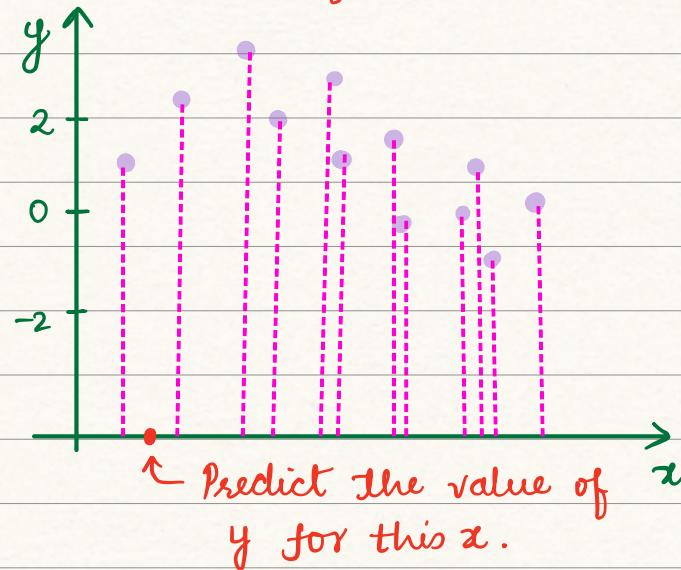
For each x_i , we are taking the value randomly lying between $[0, 1]$.

$$\sin(2\pi x_1) + n_1 = y_1 \quad n_1 \sim N(0, 1)$$

noise

$$\Rightarrow \sin(2\pi x_i) + n_i = y_i$$

This way we have artificially developed our own data set.



Non-linear \rightarrow Curve fitting Approach

$$y(x, \theta) = h(x, \theta) = \underbrace{\theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_m x^m}_{m^{\text{th}} \text{ order polynomial}}$$

But I still have no idea regarding which value of m would best fit with the training data.

$$h(x, \theta) = \sum_{j=0}^m \theta_j x^j$$

$m \rightarrow$ order of the polynomial

Subject:

1

$h(x, \theta)$ is a linear function of the coefficients and non-linear function of x .

Our goal is to find the value of θ in such a way that MSE is minimized.

Suppose I fix $m = 5$;

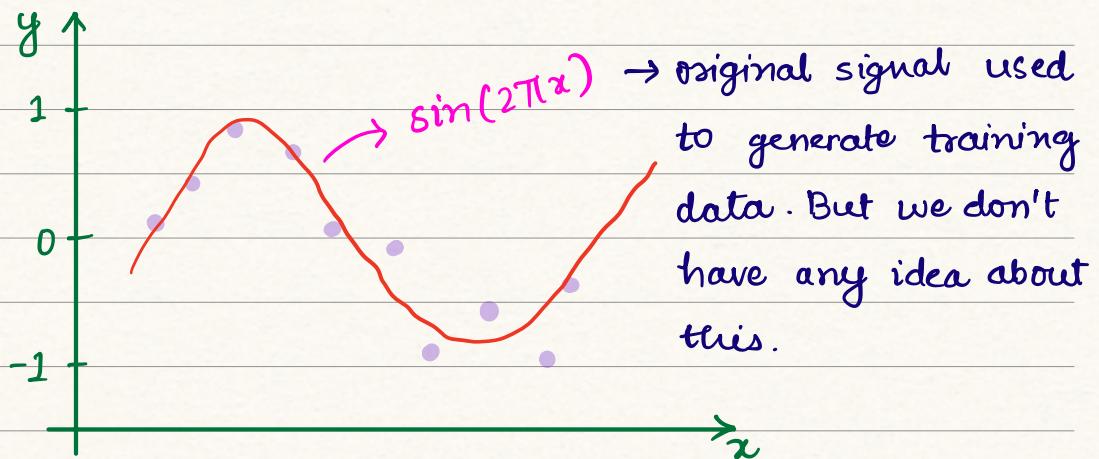
$$h(x, \theta) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \theta_5 x^5$$

$$\text{cost function, } J(\theta) = \frac{1}{2} \sum_{i=1}^n (y(x_i, \theta) - t_i)^2$$

predicted output
corresponding to
ith training example

Actual output
corresponding to
ith training example

Q.) what if we choose different values of m in the m^{th} order polynomial ?



$$m = 0 ;$$

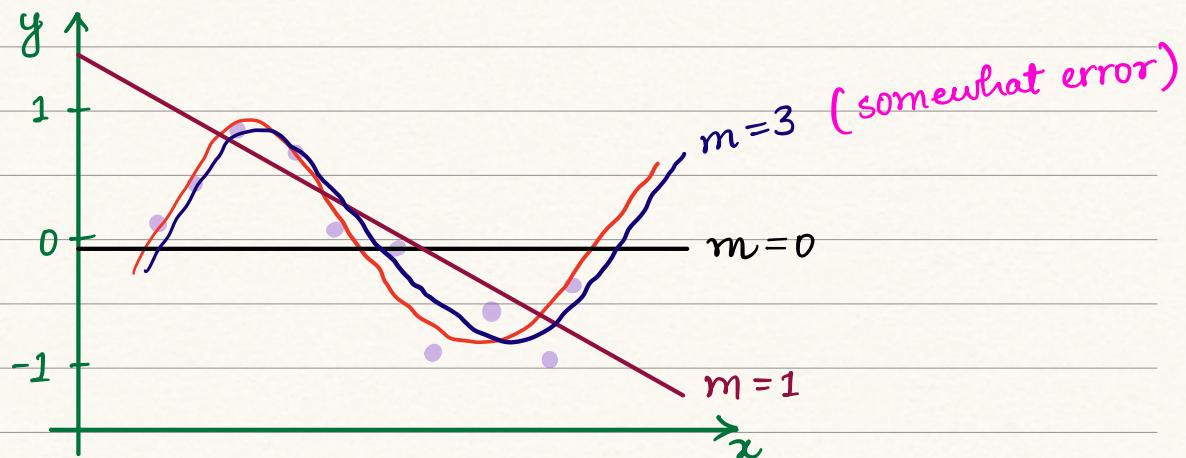
$$h(x, \theta) = \theta_0$$

$$m = 1 \text{ ;}$$

$$h(x, \theta) = \theta_0 + \theta_1 x^1$$

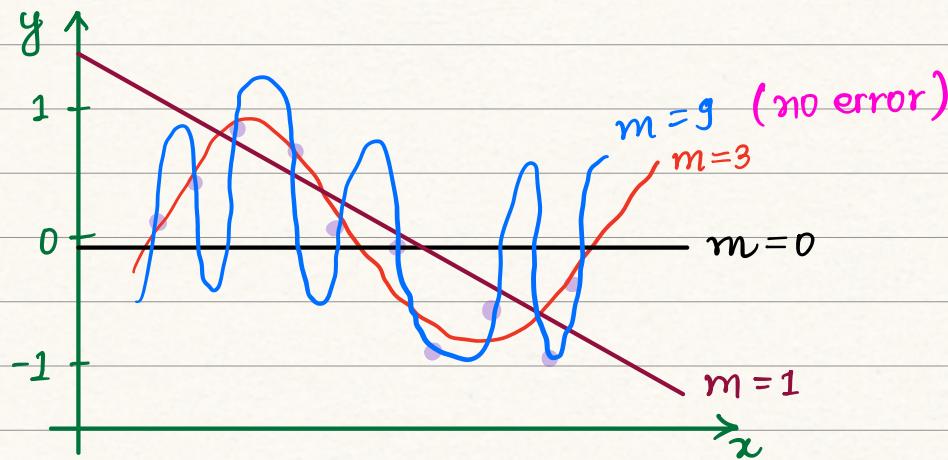
Subject:

/ /



$m = 3 ;$

$$h(x, \theta) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3$$



$m = 9 ;$

$$h(x, \theta) = \sum_{j=0}^m \theta_j x^j$$

for $m=9$; there is not any training error but there is lot of fluctuations from $\sin(2\pi x)$.

This behaviour in the context of ML is called 'overfitting'.

1000 samples (Training)

800 \rightarrow Samples I will use to train the Algorithm.

200 \rightarrow samples to test the performance of trained algorithm.

Subject:

13 /08/2025

Training Error :- $E(\theta) = \frac{1}{2N} \sum_{n=1}^N (y(x_n, \theta) - t_n)^2$

we will

compute this for $m=0, 1, 2, \dots, 9$ predicted value

$$\left. \begin{array}{l} E(\theta)_{m=0} \\ E(\theta)_{m=1} \\ \vdots \\ E(\theta)_{m=9} \end{array} \right\} \quad \text{Root mean square error}$$
$$E_{RMS} = \sqrt{\frac{2E(\theta)}{N}}$$

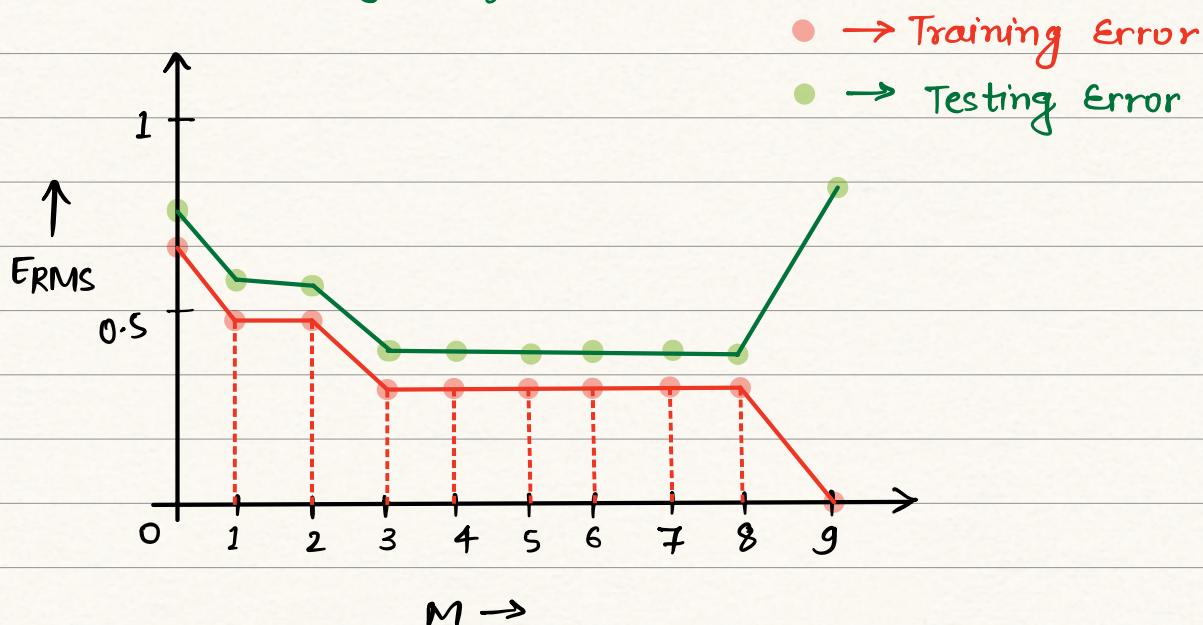
N = no. of training examples

* In this, E_{RMS} , N is included in order to compare different sizes of data sets on an equal footing and square root ensures that E_{RMS} is measured on the same scale as the target / output variable t .

Testing Error :-

$$E(\theta) = \frac{1}{2N} \sum_{n=1}^N (y(x_n, \theta) - t_n)^2$$

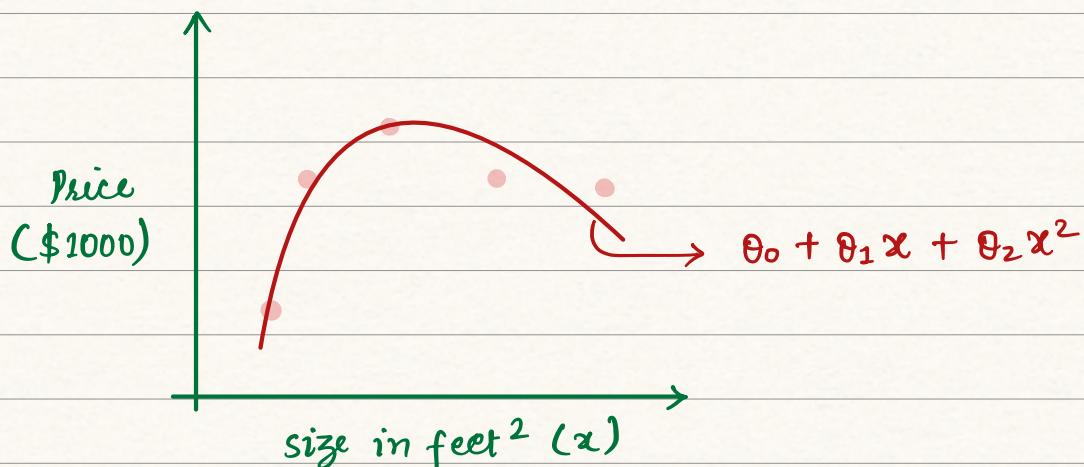
Calculate error for value of x on trained example with trained weights of θ .





This fit is very bad with the training data. This situation is called 'high bias' or 'underfit'.

underfit or high bias → Algorithm is not able to fit the training data very well.



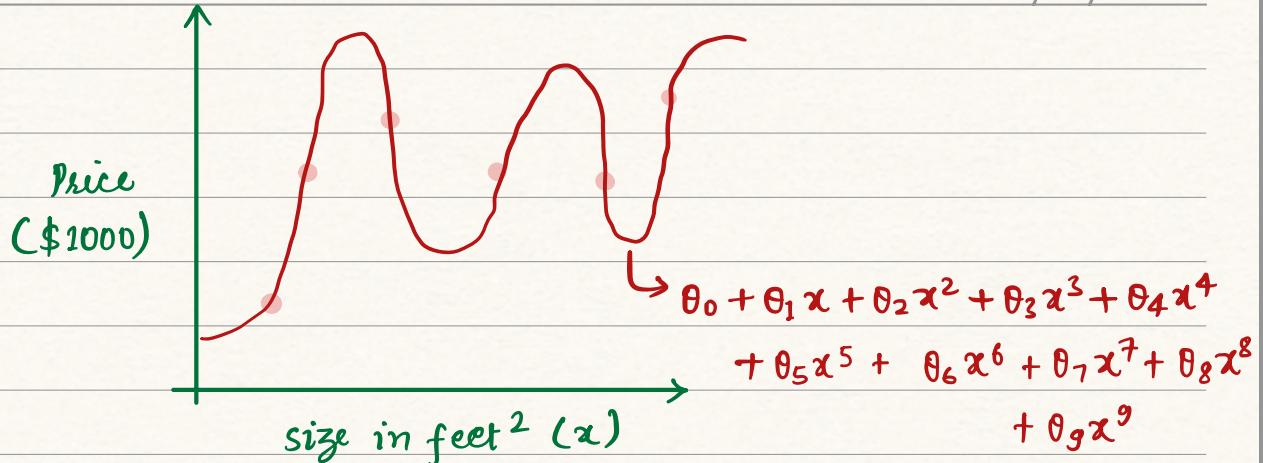
This fit with the training data is good.

→ we don't have high bias or underfit.

→ We don't have high variance or overfit.

* Our algorithm fits very well with the training data.

* Our algorithm would make very good predictions on the unseen data (or during testing phase).

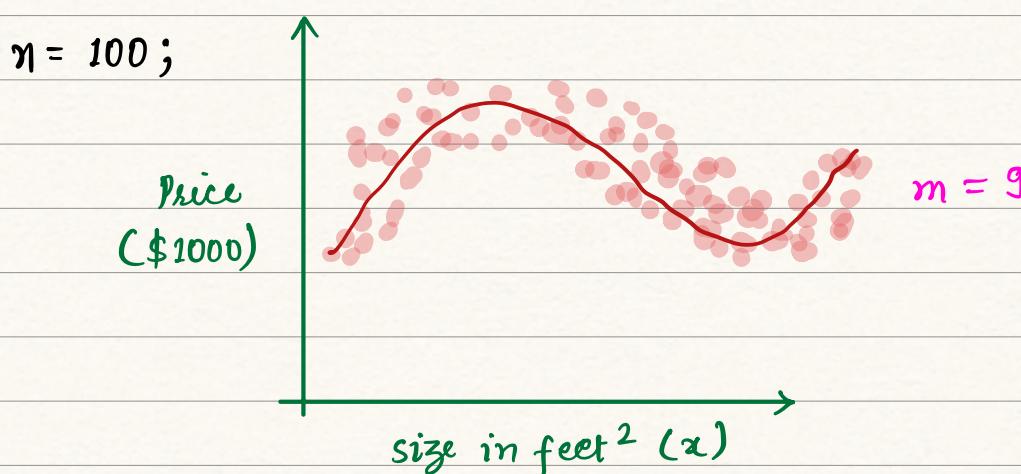
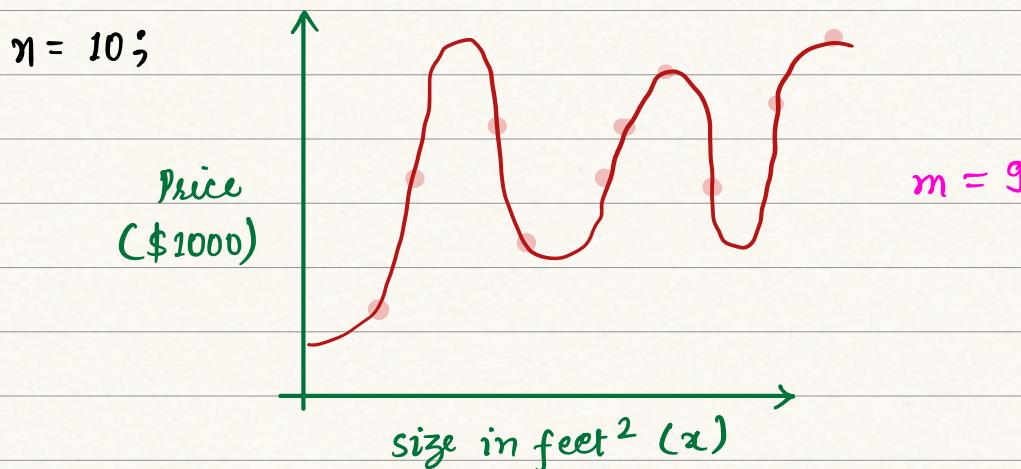


This fits the training set very well.

→ Overfitting problem \Rightarrow has high variance
 \Rightarrow leads to large testing error.

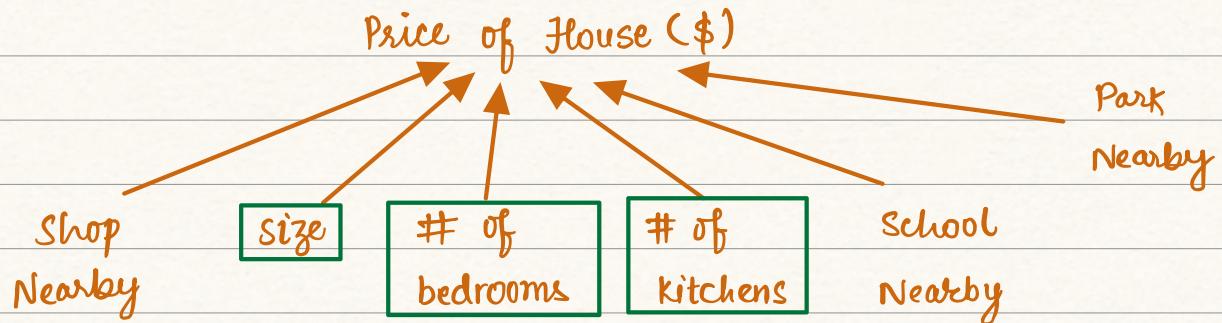
Q.) How can we overcome the overfitting problem?

Way I : Collect more training examples.

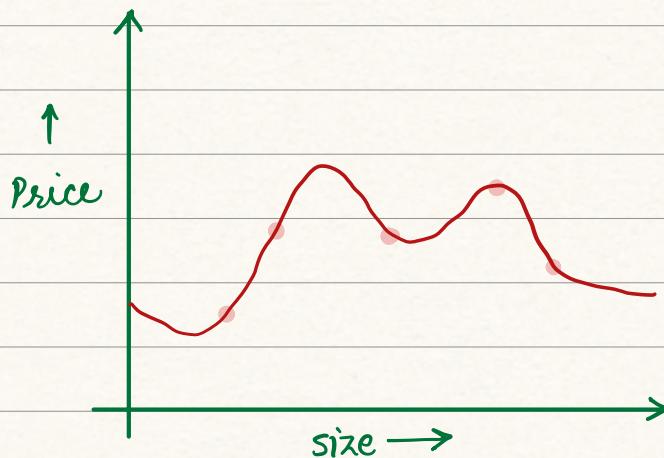


Way II: use few features

If there is a case of overfitting, only choose those features which affect the target value more than any other features.



★ Regularization —



$$h(x) = 100 + 28x - 385x^2 + 39x^3 - 174x^4$$

$$m = 4$$

Modified cost function :

$$J_{\text{Reg}}(\theta) = \frac{1}{2n} \sum_{i=1}^n (y_i - h(x_i))^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

Also called Ridge Regression.

θ_0 is
not included

λ = Regularization Parameter

$$h(x) = \theta_0 + \theta_1 x_i + \theta_2 x_i^2 + \dots + \theta_d x_i^d$$

Subject:

/ /

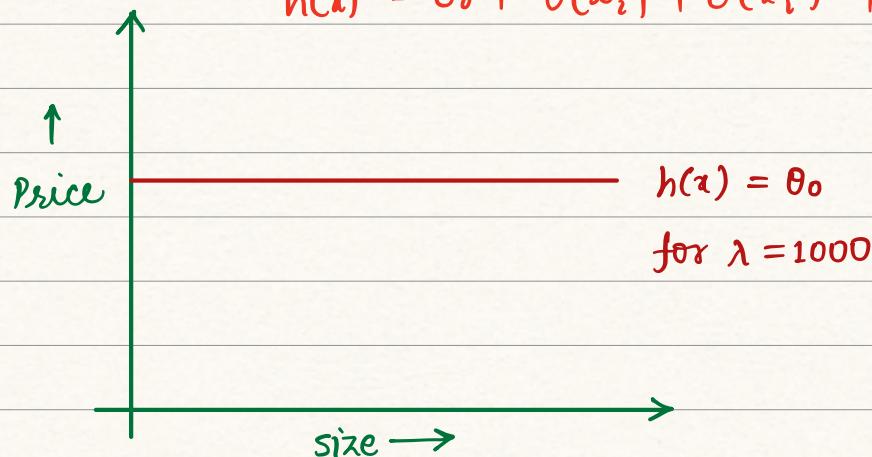
If $\lambda = 1000$;

We need to keep $\left[\frac{1000}{2} \sum_{j=1}^d \theta_j^2 \right]$ as small as possible.

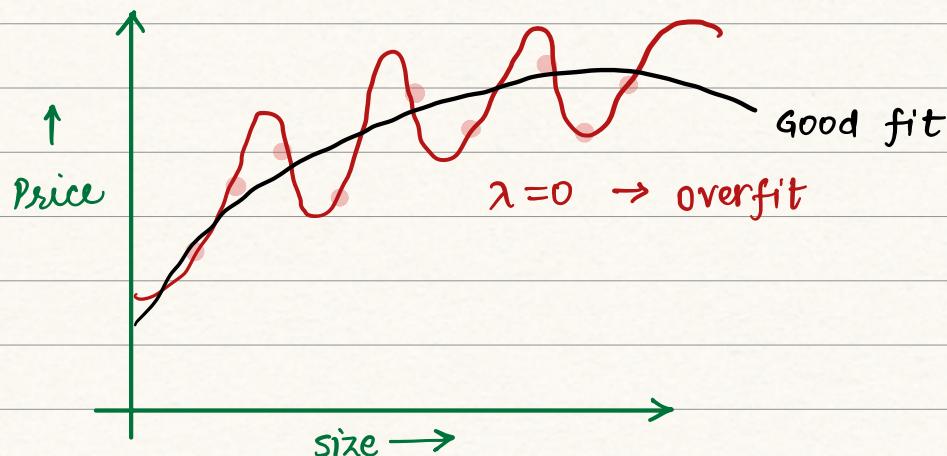
$h(x) \approx \theta_0 \rightarrow \theta_i \approx 0 \quad \forall i=1 \text{ to } d$

case of high bias

$$h(x) = \theta_0 + \theta_1(x_1) + \theta_2(x_2)^2 + \dots + \theta_d(x_d)^d$$



If $\lambda = 0$ and $d = 9$;

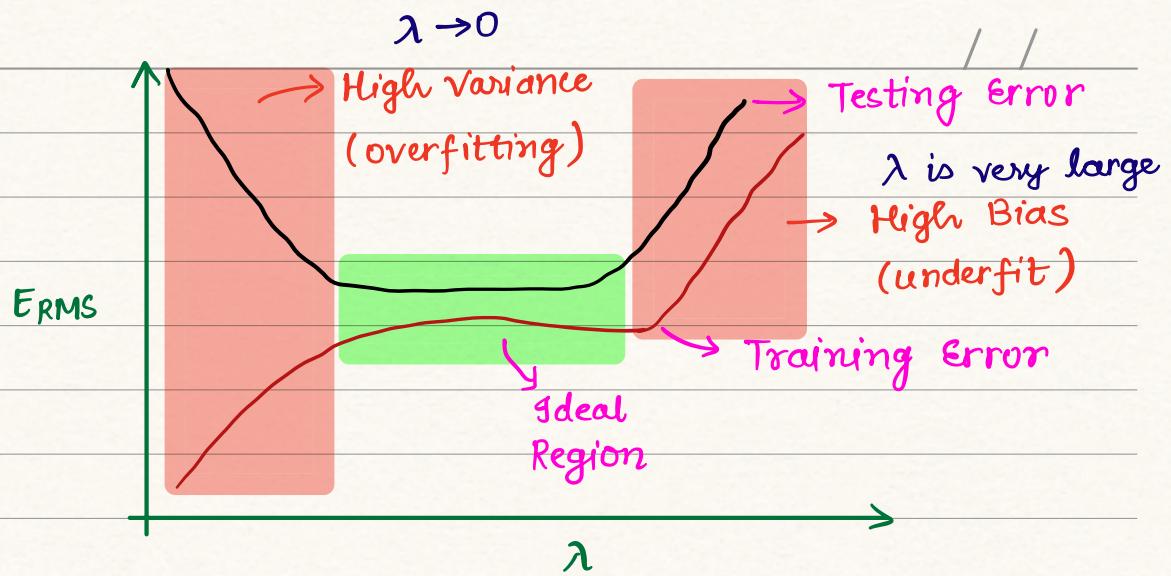


Regularized

$$\min_{\theta} J_{\text{Reg}}(\theta) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - h(x^{(i)}))^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

Regularized linear regression

Subject:



Feature Scaling (To speed up the Gradient Descent Algorithm)

Example -

Range of size of house = $300 \sim 2000 \text{ ft}^2$

Range of no. of bedrooms = $1 \sim 6$



$x_1 \rightarrow \text{size}$

$x_2 \rightarrow \text{bedrooms}$

$$x_1 (\text{scale}) = \frac{x_1}{2000 - 300}$$

$$x_2 (\text{scale}) = \frac{x_2}{6 - 1}$$

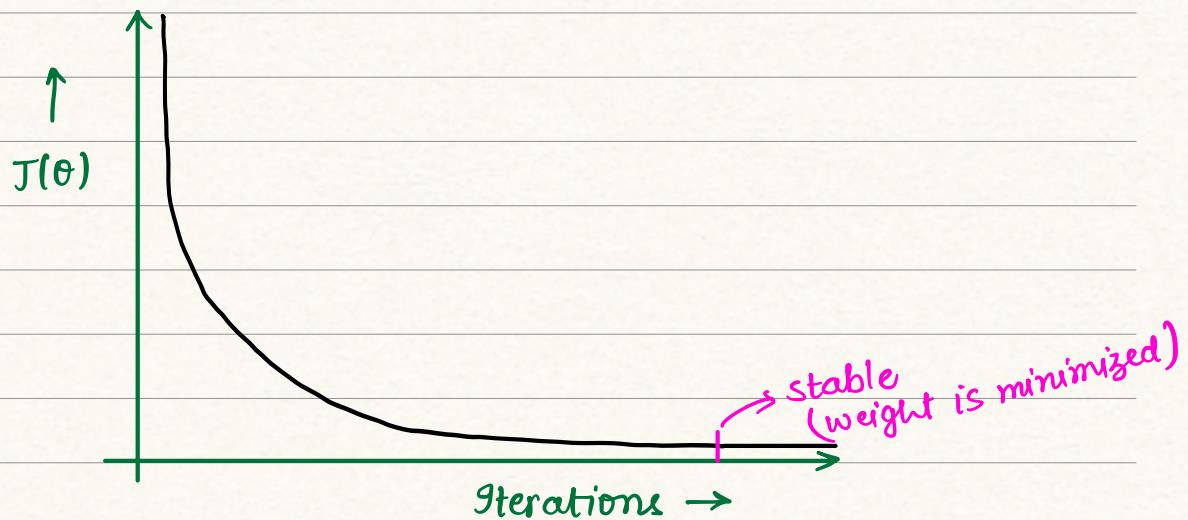
★ Mean Normalization —

μ_1 = Average for feature x_1

μ_2 = Average for feature x_2

$$x_1 \text{ (Normalized)} = \frac{x_1 - \mu_1}{x_{1\max} - x_{1\min}}$$

$$x_2 \text{ (Normalized)} = \frac{x_2 - \mu_2}{x_{2\max} - x_{2\min}}$$



★ Mean Squared Error (MSE) cost Function : —

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (y_i - h(x^{(i)}))^2$$

$$J(\theta) = \frac{1}{2n} (x\theta - y)^T (x\theta - y)$$

Theorem : A twice differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if the Hessian $\nabla^2 f$ is positive semi-definite for all $x \in \mathbb{R}^n$.

Q) What is Hessian ?

Subject:

/ /

If $f(\mathbf{x})$ has continuous second-order partial derivatives,
The Hessian of $f(\mathbf{x})$ is defined as

$$H(\mathbf{x}) = \nabla((\nabla f(\mathbf{x}))^T) \quad \nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \rightarrow \begin{array}{l} f_1 \\ f_2 \\ \vdots \\ f_n \end{array}$$

* Jacobian —

Gradient of vector valued function.

$$\mathbf{J} = \nabla((\nabla f)^T) = \nabla[f_1 \ f_2 \ \dots \ f_n]$$

$$= \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & & & \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & & & \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \frac{\partial^2 f}{\partial x_2 \partial x_n} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \begin{array}{l} \text{Hessian} \\ \text{matrix} \\ \downarrow \\ \text{symmetric} \end{array}$$

$$\nabla^T = \begin{bmatrix} \frac{\partial}{\partial x_1} & \frac{\partial}{\partial x_2} & \dots & \frac{\partial}{\partial x_n} \end{bmatrix}$$

Subject: / /

$$\nabla \nabla^T$$

$$\nabla = \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{bmatrix}$$

$$\nabla^T = \begin{bmatrix} \frac{\partial}{\partial x_1} & \frac{\partial}{\partial x_2} & \dots & \frac{\partial}{\partial x_n} \end{bmatrix}$$

$$\nabla \nabla^T = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} & \frac{\partial^2}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2}{\partial x_1 \partial x_n} \\ \frac{\partial^2}{\partial x_1 \partial x_2} & \frac{\partial^2}{\partial x_2^2} & \dots & \frac{\partial^2}{\partial x_2 \partial x_n} \\ \vdots & & & \\ \frac{\partial^2}{\partial x_1 \partial x_n} & \frac{\partial^2}{\partial x_2 \partial x_n} & \dots & \frac{\partial^2}{\partial x_n^2} \end{bmatrix}$$

Example -

$$f(x, y) = x^2y + 3y^2$$

$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 2xy \\ 6y + x^2 \end{bmatrix} \rightarrow \begin{array}{l} f_1 \\ f_2 \end{array}$$

$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{bmatrix} = \begin{bmatrix} 2y & 2x \\ 2x & 6 \end{bmatrix}$$

Subject:

/ /

* Positive Semi definite Matrix (PSD) —

If the matrix (M) is symmetric and $v^T M v \geq 0$
 $\forall v \in \mathbb{R}^n$. Then M is PSD.

$v^{n \times 1}$

$M^{n \times n}$

$v^T_{1 \times n}$

$M \rightarrow \text{Hessian Matrix}$

$$M = H(x)$$

$$v^T H(x) v \geq 0 \quad \forall v \in \mathbb{R}^n$$

Proving —

$$J(\theta) = \frac{1}{2n} (x\theta - y)^T (x\theta - y)$$

$$\nabla_{\theta}(J(\theta)) = \frac{1}{2n} (2x^T x\theta - 2x^T y)$$

$$= \frac{1}{n} (x^T x\theta - x^T y)$$

$$\nabla_{\theta}(J(\theta)) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \frac{\partial J(\theta)}{\partial \theta_1} \\ \frac{\partial J(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_d} \end{bmatrix}$$

$$y \in \mathbb{R}^n$$

$$\theta \in \mathbb{R}^{(d+1) \times 1}$$

$$x \in \mathbb{R}^{n \times (d+1)}$$

$$x^T x\theta \in \mathbb{R}^{(d+1) \times 1}$$

$$x^T y \in \mathbb{R}^{(d+1) \times 1}$$

Subject:

/ /

$$\begin{aligned}\text{Hessian of } J(\theta) &= \nabla_{\theta} (\nabla_{\theta}^T J(\theta)) \\ &= \frac{1}{n} \nabla_{\theta} (\theta^T x^T x - y^T x)\end{aligned}$$

$d=1$;

$$\theta^T x^T x = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} x_1 & x_2 \\ x_2 & x_1 \end{bmatrix}$$

$$= \begin{bmatrix} \theta_0 x_1 + \theta_1 x_2 & \theta_0 x_2 + \theta_1 x_1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial f_1}{\partial \theta_0} & \frac{\partial f_1}{\partial \theta_1} \\ \frac{\partial f_2}{\partial \theta_0} & \frac{\partial f_2}{\partial \theta_1} \end{bmatrix}$$

$$= \begin{bmatrix} x_1 & x_2 \\ x_2 & x_1 \end{bmatrix} = x^T x$$

$$\nabla^2 J(\theta) = x^T x \rightarrow \text{Hessian}$$

$v^T x^T x v$ should be ≥ 0 $\forall v \in \mathbb{R}^{(d+1)}$

Concept:

$$\text{Suppose } x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$$

L2 - Norm of vector x is

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} \geq 0$$

$$\|x\|_2^2 = x_1^2 + x_2^2 + \dots + x_n^2 \geq 0$$

$$\begin{aligned}\|x\|_2^2 &= x^T x \\ &= \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}\end{aligned}$$

$$v^T x^T x v = (x v)^T (x v)$$

$$x \in \mathbb{R}^{n \times (d+1)}$$

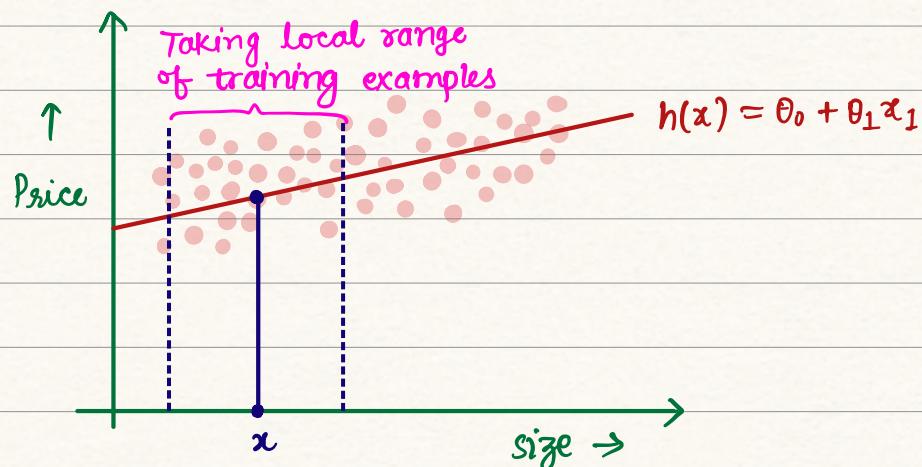
$$v \in \mathbb{R}^{(d+1) \times 1}$$

$$(x v)^T \in \mathbb{R}^{n \times 1}$$

$$\|x v\|_2^2 \geq 0$$

→ So cost function should be convex function.

* Locally Weighted Linear Regression (LWLR) —



To find the price for point x , consider only a range of local training examples.

Modified cost function :

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n w_i (y^{(i)} - h(x^{(i)}))^2$$

where $w_i = \exp\left(-\frac{(x^{(i)} - x)^2}{2\sigma^2}\right)$

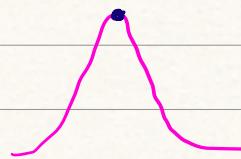
σ = Range of local training examples

σ = hyperparameter

$\sigma \rightarrow$ Defines how quickly the weight of a training example falls off with distance.

The problem here is that when we change the value of x , there will have a different range of training examples.

Decaying quickly



Decaying slowly



Needs more no. of training examples

Needs less no. of training examples

Locally weighted Linear Regression is also called non-parametric algorithm.

Unweighted (conventional) Linear Regression is called parametric Learning Algorithm.

* Logistic Regression —

for classification

Example :-

Email spam/not spam

Online Transactions - Fraudulent (Yes/No)

Disease - Tumor (Yes/No)

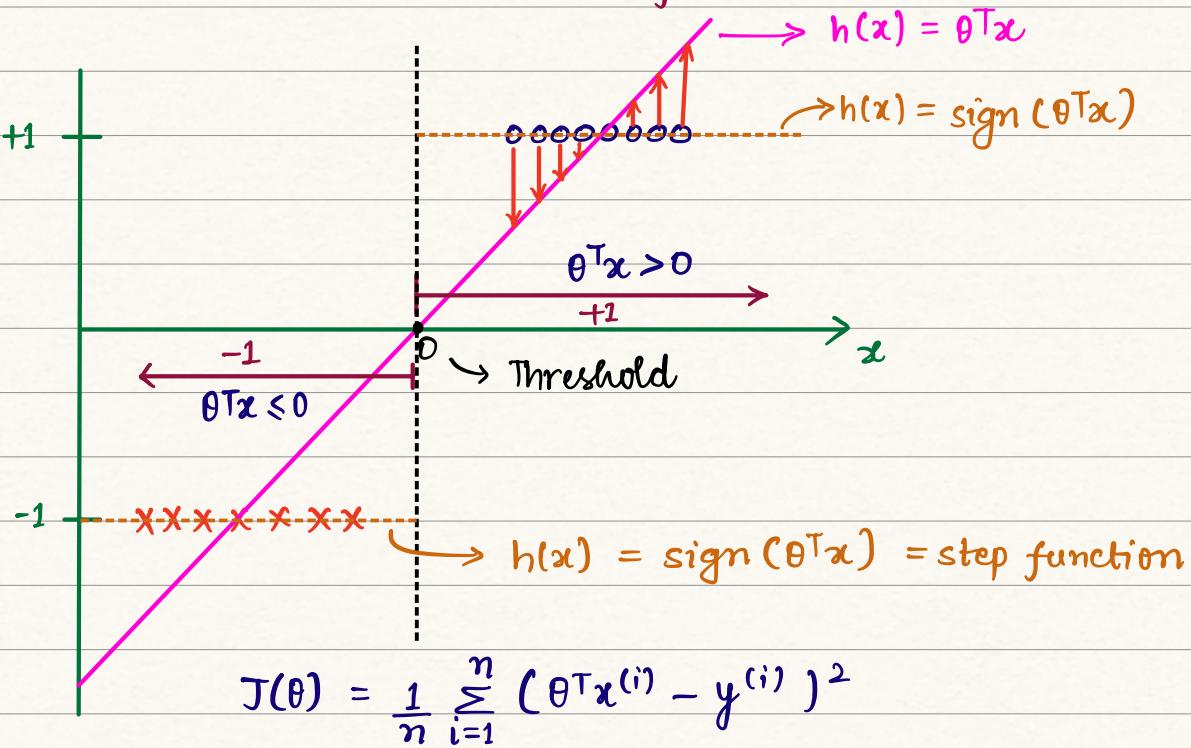
Features		Output
x_1	x_2	y (Disease)
a	b	Yes
c	d	No
e	f	Yes
g	h	Yes

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

$0 \rightarrow$ Negative class
 $1 \rightarrow$ Positive class

$y \in \{0, 1, 2, \dots\}$

Multi-class classification

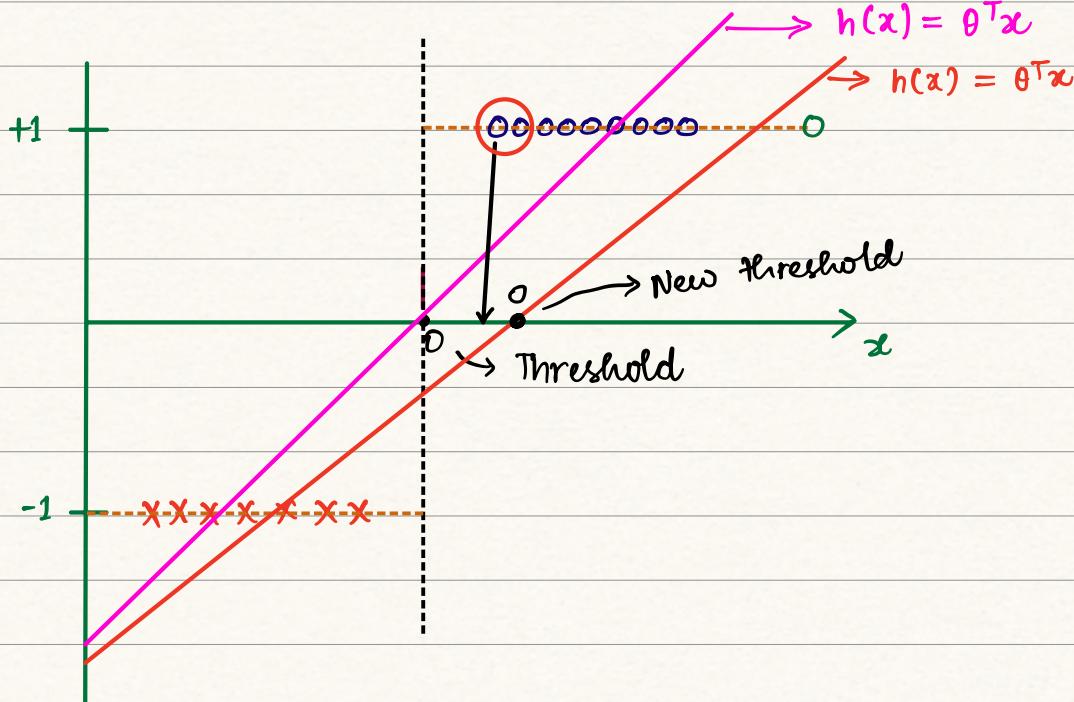


Subject:

$$\text{Prediction, } h(\mathbf{x}) = \text{sign}(\theta^T \mathbf{x}) = \begin{cases} +1 & ; \theta^T \mathbf{x} > 0 \\ -1 & ; \theta^T \mathbf{x} \leq 0 \end{cases}$$

Q) what is the problem with this method of classification?

→ Sensitive to outlier



$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$

$$h(\mathbf{x}^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

For testing the house price based on the size of house, we were taking $h(\mathbf{x}) = \theta^T \mathbf{x}$ for predicting the house price (we were taking continuous range of values).

$$h(\mathbf{x}) = \theta^T \mathbf{x}$$

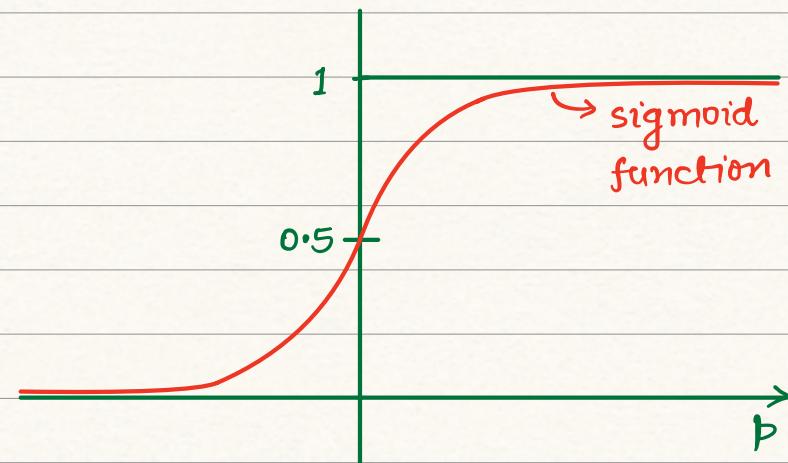
But for logistic regression, we want

$$0 \leq h(\mathbf{x}) \leq 1$$

because $y \in \{0, 1\}$ by doing $h(\mathbf{x}) = g(\theta^T \mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$

Subject:

/ /



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

$$p = \theta^T x$$

$$h(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}} = P(y=1|x; \theta)$$

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

Probability that $y=1$
given x & parameterized
by θ

Predict " $y=1$ " if $h(x) = \frac{1}{1+e^{-\theta^T x}} \geq 0.5$

Predict " $y=0$ " if $h(x) = \frac{1}{1+e^{-\theta^T x}} < 0.5$

Testing ;

$$g(h(x)) = \frac{1}{1+e^{-\theta^T x}} \geq 0.5 \Rightarrow x \text{ belongs to class 1}$$

$$g(h(x)) = \frac{1}{1+e^{-\theta^T x}} < 0.5 \Rightarrow x \text{ belongs to class 0}$$

$$x = \begin{bmatrix} 1 \\ x_1 \end{bmatrix}$$

x_1 = Tumor size

$$g(h(x)) = \frac{1}{1+e^{-\theta^T x}} = 0.7 \quad \{ \text{The person has a tumor?} \}$$

Subject:

26/08/2025

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1+e^{-\theta^T x^{(i)}}} - y^{(i)} \right)^2$$

Problem here is this function is no longer a convex function.

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

$$x^{(i)} \in \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad x_0^{(i)} = 1$$

$\forall i \in \{1, 2, 3, \dots, m\}$

$$y^{(i)} \in \{0, 1\} ; i \in \{1, 2, 3, \dots, m\}$$

$$h(x) = \frac{1}{1 + e^{-\theta^T x}} \rightarrow \text{This is used for prediction purpose.}$$

To find the optimum value of θ , we need to first find out the cost function which should be minimized w.r.t θ so as to find the optimum value of θ .

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

$$\frac{dg(s)}{ds} = \frac{d}{ds} \left(\frac{1}{1 + e^{-s}} \right)$$

$$= \frac{e^{-s}}{(1 + e^{-s})^2}$$

$$= \frac{1}{(1 + e^{-s})} \left[1 - \frac{1}{1 + e^{-s}} \right]$$

$$g'(s) = g(s)[1 - g(s)]$$

Subject:

/ /

Data set for logistic Regression

1st training example	$\left\{ x_0^{(1)} \quad x_1^{(1)} \quad x_2^{(1)} \dots x_n^{(1)} \right.$	$y^{(1)}$	output values
	$x_0^{(2)} \quad x_1^{(2)} \quad x_2^{(2)} \dots x_n^{(2)}$	$y^{(2)}$	$y^{(i)} \in \{0,1\}$
	\vdots	\vdots	
	$x_0^{(m)} \quad x_1^{(m)} \quad x_2^{(m)} \dots x_n^{(m)}$	$y^{(m)}$	

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

$$P(y=0 | x; \theta) = 1 - \left[\frac{1}{1+e^{-\theta^T x}} \right]$$

output $y \in \{0,1\}$ is a Bernoulli Distributed Random variable.

Basics of Probability :-

If X is a Bernoulli distributed Random variable with $P(X=1) = p$, $P(X=0) = 1-p = q$. Then the probability mass function of Bernoulli distribution is

$$f(k, p) = p^k q^{(1-k)}$$

$$= p^k (1-p)^{(1-k)} \quad k \in \{0, 1\}$$

$$f(y | x; \theta) = (h(x))^y (1-h(x))^{1-y}$$

$$y \in \{0, 1\}$$

x, y are independent

Subject:

/ /

$$f_{x,y}(x, y) = f_x(x) f_y(y)$$

$$f(y|x; \theta) = \prod_{i=1}^m (h(x^{(i)}))^{y^{(i)}} (1-h(x^{(i)}))^{1-y^{(i)}}$$

considering all y

$$= L(\theta | x; y)$$

log likelihood ;

$$\ell(\theta) = \log(L(\theta | x; y))$$

$$= \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1-y^{(i)}) \log (1-h(x^{(i)}))$$

To maximize the log likelihood, we would consider the Gradient ascent

update Rule :

$$\theta := \theta + \alpha \nabla_{\theta} \ell(\theta)$$

for simplification, suppose $m=1$ (i.e., only one training example)

$$\ell(\theta) = y \log(h(x)) + (1-y) \log(1-h(x))$$

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = \frac{y}{h(x)} \frac{\partial}{\partial \theta_j} (h(x)) + \frac{(1-y)}{1-h(x)} \cdot (-1) \cdot \frac{\partial}{\partial \theta_j} h(x)$$

$$\frac{\partial h(x)}{\partial \theta_j} = h(x)(1-h(x)) \cdot \frac{\partial}{\partial \theta_j} (\theta^T x)$$

$$= h(x)(1-h(x)) \cdot x_j \quad \begin{cases} \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ \frac{\partial}{\partial \theta_2} = x_2 \end{cases}$$

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = \left[\frac{y}{h(x)} - \frac{(1-y)}{1-h(x)} \right] (h(x)(1-h(x)) \cdot x_j)$$

Subject:

/ /

$$\Rightarrow \theta_j := \theta_j + \alpha (y - h(x)) \cdot x_j$$

$$J(\theta) = \log(L(\theta))$$

$$= \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1-y^{(i)}) \log (1-h(x^{(i)}))$$

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

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

Logistic loss, Loss = $-J(\theta)$

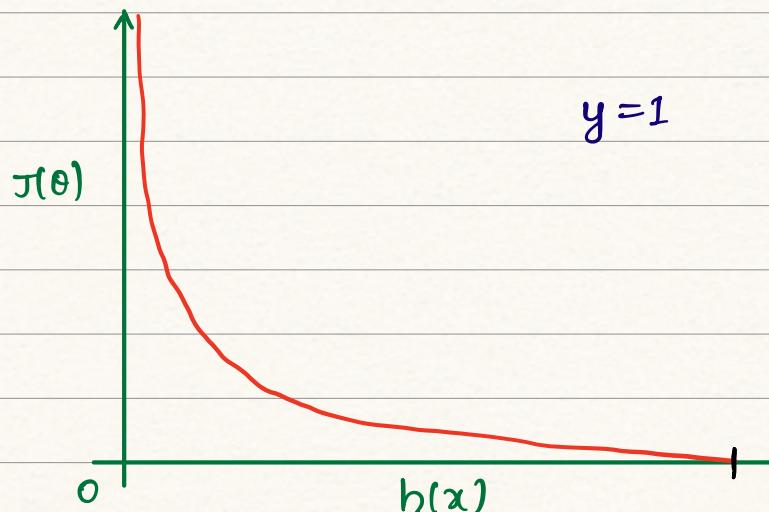
↳ Binary cross entropy
loss function

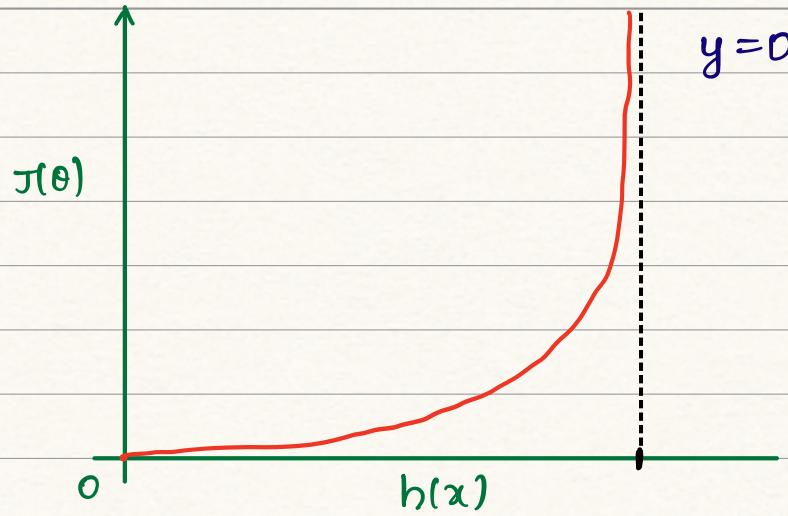
for m=1;

$$J(\theta) = -J(\theta) = -y \log(h(x)) - (1-y) \log(1-h(x))$$

$$y = 1 \longrightarrow J(\theta) = -\log(h(x))$$

$$y = 0 \longrightarrow J(\theta) = -\log(1-h(x))$$





* Binary Cross Entropy Loss function : —

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

m = no. of training examples

$y^{(i)}$ = True label corresponding to the training example

$h(x^{(i)})$ = Predicted probability corresponding to i th training example

$$\begin{aligned} J(\theta) &= - \left(\sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1-y^{(i)}) \log (1-h(x^{(i)})) \right) \\ &= - \left(\sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + \log (1-h(x^{(i)})) - y^{(i)} \log (1-h(x^{(i)})) \right) \\ &= - \left(\underbrace{\sum_{i=1}^m y^{(i)} \log \frac{h(x^{(i)})}{1-h(x^{(i)})}}_{I} + \underbrace{\log (1-h(x^{(i)}))}_{II} \right) \end{aligned}$$

Solving I ;

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

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

Subject:

$$\log \left(\frac{\frac{1}{1+e^{-\theta^T x^{(i)}}}}{1 - \frac{1}{1+e^{-\theta^T x^{(i)}}}} \right) \equiv \log(e^{+\theta^T x^{(i)}}) \\ \Rightarrow \theta^T x^{(i)}$$

$$\theta^T x^{(i)} = \theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_d x_d^{(i)}$$

$$\nabla_{\theta}(x) = \begin{bmatrix} \frac{\partial x}{\partial \theta_0} \\ \frac{\partial x}{\partial \theta_1} \\ \vdots \\ \frac{\partial x}{\partial \theta_d} \end{bmatrix} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_d^{(i)} \end{bmatrix} = x^{(i)}$$

$$v^T H v \geq 0$$

For function to be convex

Hessian is positive semi definite so function is convex
for $\forall v \in \mathbb{R}^{(d+1)}$.

$\Rightarrow \theta^T x^{(i)}$ is convex

$\rightarrow (-ve)$ comes from the whole expression.
 $F = -\log(1-h(x^{(i)}))$ is convex.

I want to prove this. How can we prove this?

$$\begin{aligned} F &= -\log(1-h(x^{(i)})) \\ &= -\log\left(1 - \frac{1}{1+e^{-\theta^T x^{(i)}}}\right) \\ &= -\log\left(\frac{e^{-\theta^T x^{(i)}}}{1+e^{-\theta^T x^{(i)}}}\right) \\ &= -\log(e^{-\theta^T x^{(i)}}) - \log(1+e^{-\theta^T x^{(i)}}) \end{aligned}$$

Subject:

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

$$\nabla_{\theta}(F) = x^{(i)} + \nabla_{\theta} \left(\log \left(1 + e^{-\theta^T x^{(i)}} \right) \right)$$

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

$$= x^{(i)} \left[1 - \frac{e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}} \right]$$

$$= \boxed{\frac{1}{1 + e^{-\theta^T x^{(i)}}}} \cdot x^{(i)} \rightarrow \text{scalar}$$

$$= h(x^{(i)}) \cdot x^{(i)}$$

Hessian is $\nabla_{\theta}(\nabla_{\theta}^T(F))$

$$\nabla_{\theta}(F) = h(x^{(i)}) \cdot x^{(i)}$$

$$\nabla_{\theta}^T(F) = h(x^{(i)}) \cdot (x^{(i)})^T$$

$$\nabla_{\theta}(\nabla_{\theta}^T(F)) = \nabla_{\theta}(h(x^{(i)})) \cdot x^{(i)T}$$

$$= \nabla_{\theta} \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} \right) \cdot x^{(i)T}$$

$$= \frac{e^{-\theta^T x^{(i)}}}{(1 + e^{-\theta^T x^{(i)}})^2} \cdot x^{(i)} \cdot x^{(i)T}$$

$$\equiv \left(\frac{1}{e^{-\theta^T x^{(i)}}} \right) \left(1 - \frac{1}{1 + e^{-\theta^T x^{(i)}}} \right) \cdot x^{(i)} \cdot x^{(i)T}$$

Subject:

$$= h(x^{(i)}) (1 - h(x^{(i)})) \cdot x^{(i)} \cdot x^{(i)T}$$

↓
Scalar

$$x^{(i)} \in \mathbb{R}^{(d+1) \times 1}$$

$$x^{(i)T} \in \mathbb{R}^{1 \times (d+1)}$$

$x^{(i)} \cdot x^{(i)T}$ is symmetric

$$= h(x^{(i)}) (1 - h(x^{(i)})) \cdot x^{(i)} \cdot x^{(i)T}$$

$$= p$$

Next task is to prove that $v^T p v \geq 0$; $v \in \mathbb{R}^{(d+1) \times 1}$

$$v^T x^{(i)} \cdot x^{(i)T} v \equiv \|v^T x^{(i)}\|_2^2 \geq 0 \quad \{v^T p v = \text{scalar}\}$$

So F is convex.

$$= - \left(\underbrace{\sum_{i=1}^m y^{(i)} \log \frac{h(x^{(i)})}{1-h(x^{(i)})}}_{\text{Convex}} + \underbrace{\log(1-h(x^{(i)}))}_{\text{Convex}} \right)$$

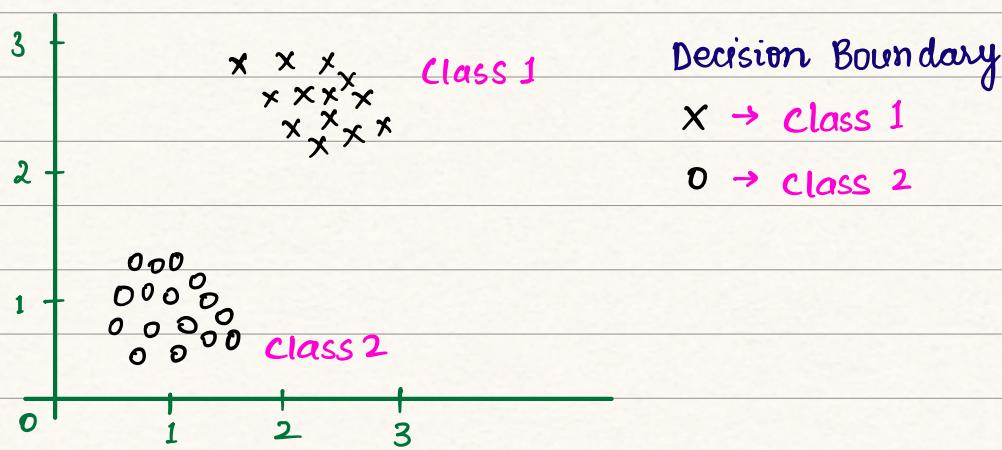
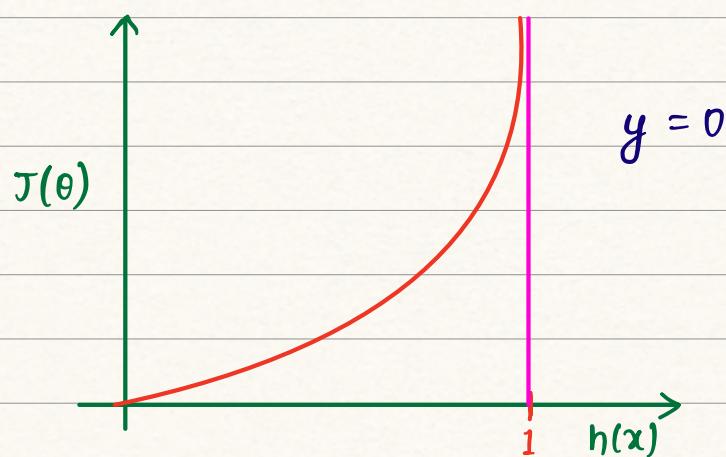
$$x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_d^{(i)} \end{bmatrix}$$

$$\|x^{(i)}\|_2^2 = 1 + (x_1^{(i)})^2 + (x_2^{(i)})^2 + \dots + (x_d^{(i)})^2$$

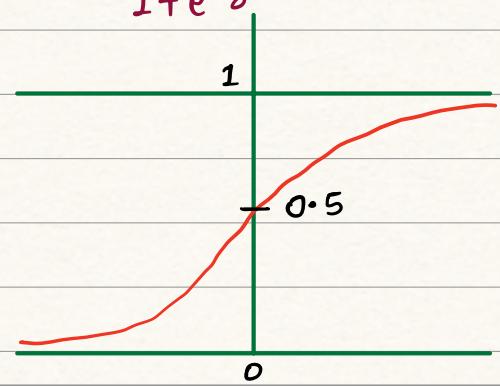
$$J(\theta) = -y \log h(x) - (1-y) \log(1-h(x))$$

$$J(\theta) = \begin{cases} -\log h(x) ; \text{ if } y=1 \\ -\log(1-h(x)) ; \text{ if } y=0 \end{cases}$$

Subject: / /



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



Subject:

/ /

Predict '1' when $g(z) \geq 0.5$

Predict '0' when $g(z) < 0.5$

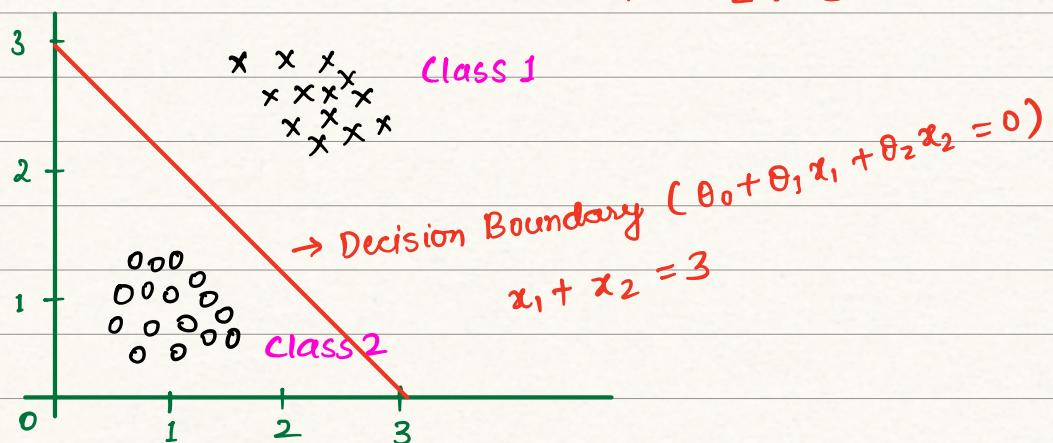
$$h(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$= \frac{1}{1 + e^{-\theta^T x}} = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}$$

Suppose you get $\theta_0 = -3, \theta_1 = 1, \theta_2 = 1$

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

Predict '1' when
 $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$
 $-3 + 1x_1 + 1x_2 \geq 0$
 $x_1 + x_2 \geq 3$



Testing

Suppose you get $x^{(i)}$

$$\text{Then } h_{\theta}(x^{(i)}) = \frac{1}{1 + e^{-\theta^T x^{(i)}}} \geq 0.5 \rightarrow y = 1$$

we know θ now

after training

$$x_1 + x_2 = 3$$

$$x_1^{(i)} + x_2^{(i)} > 3$$

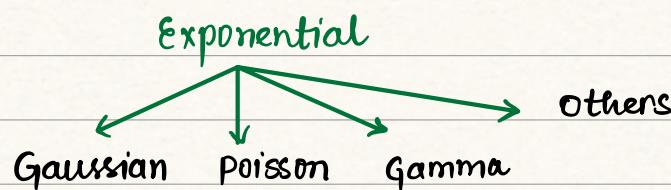
★ Exponential Family —

* Distribution falls under the exponential family if it can be written of the form :—

$$f(y; \boldsymbol{\alpha}) = b(y) \exp(\eta^T T(y) - a(\eta))$$

$\eta \rightarrow$ Natural parameter

$T(y) = y$ (sufficient statistic)



Q.) Gaussian and Bernoulli are the family of exponential distribution.
Prove it.

Bernoulli Distribution

$$\begin{aligned} y &\in \{0, 1\} \\ p(y=1 | \boldsymbol{x}; \theta) &= p(y=1) = \phi \\ p(y=0 | \boldsymbol{x}; \theta) &= 1 - p(y=1) = 1 - \phi \end{aligned}$$

$$f(y; \phi) = \phi^y (1-\phi)^{1-y} ; y \in \{0, 1\}$$

We want this expression to be in the form of

$$f(y; \boldsymbol{\alpha}) = b(y) \exp(\eta^T T(y) - a(\eta))$$

$$\begin{aligned} f(y; \phi) &= \phi^y (1-\phi)^{1-y} \\ &= \exp[\log(\phi^y (1-\phi)^{1-y})] \\ &= \exp[y \log \phi + (1-y) \log(1-\phi)] \\ &= \exp[y \log \phi - y \log(1-\phi) + \log(1-\phi)] \\ &= \exp\left[y \log\left(\frac{\phi}{1-\phi}\right) + \log(1-\phi)\right] \end{aligned}$$

$$b(y) = 1$$

$$\eta = \log\left(\frac{\phi}{1-\phi}\right)$$

Subject:

/ /

$$T(y) = y \quad a(\eta) = \log(1-\phi)$$

$$\phi = \frac{1}{1+e^{-x}} \quad \left. \right\} p(y=1)$$

sigmoid function

Gaussian Distribution ($\sigma^2 = 1$)

$$\begin{aligned} f(y; \mu, \sigma^2) &= f(y; \mu) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y^2 + \mu^2 - 2y\mu)}{2}\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) \exp\left(\mu y - \frac{\mu^2}{2}\right) \end{aligned}$$

$$b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$$

$$T(y) = y$$

$$\eta = \mu$$

$$a(\eta) = \frac{\mu^2}{2} = \frac{\eta^2}{2}$$

* Generalized Linear Models —

$$h(x) = \theta^T x \rightarrow \text{Linear Regression}$$

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

Belong to the class of Generalized linear models.

Subject:

/ /

In Regression ;

$$y|x; \theta \sim N(\mu, \sigma^2)$$
$$\mu = \theta^T x$$

In logistic ;

$$y|x; \theta \sim \text{Bernoulli}(p)$$

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

Constructing Generalized Linear Models —

* Following assumptions are made :

1.) $y|x; \theta \sim \text{Exponential family } (\eta)$

2.) For a given x , prediction $h(x) = E[y|x]$

Discrete \leftarrow

$$x = \{x_1, x_2\}$$
$$\downarrow \quad \downarrow$$
$$p_1 \quad p_2$$

Random Variable

$$E[x] = p_1 x_1 + p_2 x_2$$

$$E[y] = 1 \cdot \phi + 0 \cdot (1-\phi) \quad y \in \{0, 1\}$$

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

3.) The natural parameter η and the input vector x is related as $\eta = \theta^T x$.

* Linear Regression is a special case of Generalized Linear models.

1.) $y|x; \theta \sim \text{Exponential } (\eta)$

$\sim \text{Gaussian } (\mu)$

$$\Rightarrow \mu = \eta$$

Subject:

/ /

$$2.) h(x) = E[y|x] = u = \eta$$

$$3.) \eta = \theta^T x$$

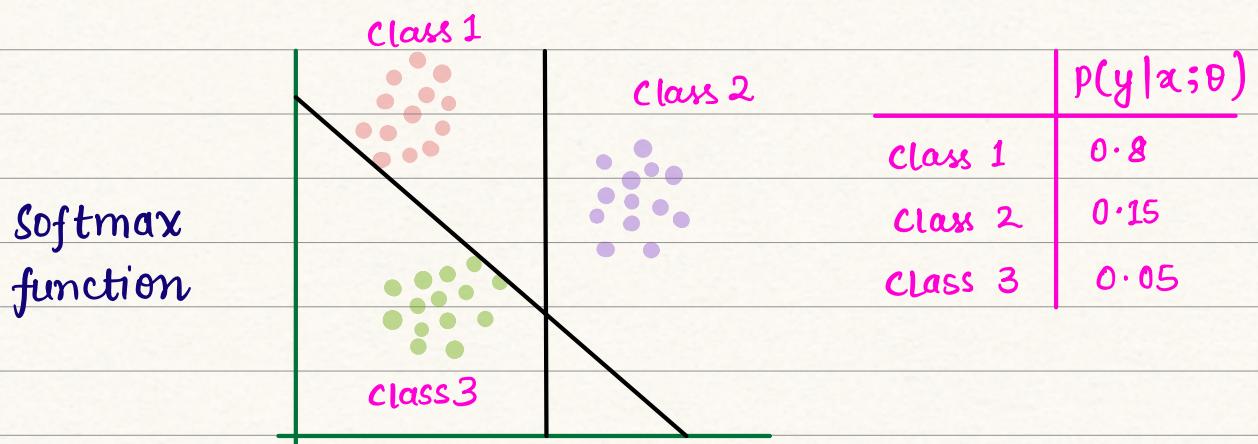
* Logistic Regression —

$$1.) y|x; \theta \sim \text{Exponential}(\eta)$$
$$\sim \text{Bernoulli}(\phi)$$

$$2.) h(x) = E[y|x] = \phi = \frac{1}{1+e^{-\eta}}$$

$$3.) \eta = \theta^T x$$
$$\Rightarrow h(x) = \frac{1}{1+e^{-\theta^T x}}$$

* Multi-class classification —



* Softmax Regression —

Multi-class classification

$y \in \{1, 2, \dots, k\} \rightarrow \text{Discrete Value}$

Example :-

This can be modeled with
Multinomial Distribution.



Task 1 : Let's represent multinomial distribution as a class of exponential distribution.

* Multinomial distribution over k possible outcomes.

Let $\phi_1, \phi_2, \dots, \phi_k$

$$P(y=1) \quad P(y=2) \quad P(y=k)$$

$$\phi_i = P(y=i ; \phi)$$

$$\phi_k = P(y=k ; \phi)$$

$$\left\{ \sum_{i=1}^k \phi_i = 1 \right.$$

$$\Rightarrow \phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$$

Exponential Distribution :

$$f(y ; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

↑ ↑
Vector vector

To express the multinomial as an Exponential family distribution, we will define ;

$$T(y) \in \mathbb{R}^{(k-1)}$$

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \text{First Entry}$$

$$T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \text{Second Entry}$$

$$T(k-1) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \rightarrow (k-1)^{\text{th}} \text{ entry}$$

Subject:

/ /

$T(y)_i \rightarrow$ denotes the i th element of vector $T(y)$.

$$T(y)_i = \underbrace{1\{y=i\}}_{\text{Indicator function}}$$

Indicator Function : —

An indicator function $1\{\cdot\}$ takes on the value 1 if its argument is true & 0 otherwise.

$$1\{\text{True}\} = 1$$

$$1\{\text{False}\} = 0$$

$$1\{2=3\} = 0$$

$$1\{3=3\} = 1$$

$$1\{y, i\} = \begin{cases} 1 & \text{if } y=i \\ 0 & \text{if } y \neq i \end{cases}$$

$$\begin{aligned} E[1\{y, i\}] &= 1 \cdot P(y=i) + 0 \cdot P(y \neq i) \\ &= P(y=i) \\ &= \phi_i \end{aligned}$$

$$f(y, \phi) = \phi_1^{1\{y=1\}} \cdot \phi_2^{1\{y=2\}} \cdot \phi_3^{1\{y=3\}} \dots \phi_K^{1\{y=k\}}$$

For $K=2$:

$$f(y, \phi) = \phi^{1\{y=1\}} \cdot \phi^{1\{y=2\}}$$

$$\begin{aligned} 1\{y=k\} &= 1 - \sum_{i=1}^{k-1} 1\{y=i\} \\ &= \phi_1^{1\{y=1\}} \cdot \phi_2^{1\{y=2\}} \dots \phi_K^{1 - \sum_{i=1}^{k-1} 1\{y=i\}} \end{aligned}$$

$$= \phi_1^{(T(y))_1} \cdot \phi_2^{(T(y))_2} \dots \phi_K^{1 - \sum_{i=1}^{K-1} (T(y))_i}$$

$$= \exp[(T(y))_1 \log \phi_1 + (T(y))_2 \log \phi_2 + \dots + (T(y))_{K-1} \log \phi_{K-1}]$$

$$+ \log \phi_K - \log \phi_K \sum_{i=1}^{K-1} (T(y))_i]$$

Subject:

$$= \exp \left[(\tau(y))_1 \log \frac{\phi_1}{\phi_K} + (\tau(y))_2 \log \frac{\phi_2}{\phi_K} + \dots + (\tau(y))_{K-1} \log \frac{\phi_{K-1}}{\phi_K} + \log \phi_K \right]$$

$$= b(y) \exp(\eta^\top \tau(y) - a(\eta))$$

$$\eta = \begin{bmatrix} \log \left(\frac{\phi_1}{\phi_K} \right) \\ \log \left(\frac{\phi_2}{\phi_K} \right) \\ \vdots \\ \log \left(\frac{\phi_{K-1}}{\phi_K} \right) \end{bmatrix}$$

$$\tau(y) = \begin{bmatrix} (\tau(y))_1 \\ (\tau(y))_2 \\ \vdots \\ (\tau(y))_{K-1} \end{bmatrix}$$

$$\eta_i = \log \frac{\phi_i}{\phi_K} \quad i = 1, 2, \dots, K-1$$

$$a(\eta) = -\log(\phi_K)$$

$$\eta_K = \log \frac{\phi_K}{\phi_K}$$

$$b(y) = 1$$

$$\phi_i = \phi_K \exp(\eta_i)$$

$$\Rightarrow \phi_K = \frac{\phi_i}{\exp(\eta_i)}$$

$$\sum_{i=1}^K \phi_i = 1 \quad \Rightarrow \quad \phi_K \sum_{i=1}^K \exp(\eta_i) = 1$$

$$\phi_K = \frac{1}{\sum_{i=1}^K \exp(\eta_i)}$$

Subject:

$$\phi_i = \frac{\exp(\eta_i)}{\sum_{i=1}^k \exp(\eta_i)} \quad \left. \right\} \text{softmax Function}$$

$$\eta_i = \theta_i^T x \quad ; \quad i = 1, 2, \dots, k-1$$

$$\eta_k = \theta_k^T x = 0$$

$$\Rightarrow \theta_k = 0 \quad (\text{zero vector})$$

$$P(y=i | x; \theta) = \phi_i = \frac{\exp(\eta_i)}{\sum_{j=1}^k \exp(\eta_j)} = \frac{\exp(\theta_i^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)}$$

$$h(x) = E[T(y) | x; \theta]$$

$$E \begin{bmatrix} 1 \{y=1\} \\ 1 \{y=2\} \\ \vdots \\ 1 \{y=k-1\} \end{bmatrix} \mid x; \theta \Rightarrow \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

$$\begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix} = \frac{1}{\sum_{j=1}^k \exp(\theta_j^T x)} \begin{bmatrix} \exp(\theta_1^T x) \\ \exp(\theta_2^T x) \\ \vdots \\ \exp(\theta_{k-1}^T x) \end{bmatrix}$$

$$P(y=k | x; \theta) = 1 - \sum_{i=1}^{k-1} P(y=i | x; \theta)$$

$$P(y=1 | x; \theta) = \frac{\exp(\theta^T x)}{1 + \exp(\theta_1^T x)} = \frac{1}{1 + \exp(-\theta_1^T x)} = h(x)$$

$$\exp(0) = 1$$

$$P(y=0 | x; \theta) = \frac{1}{1 + \exp(\theta^T x)} = 1 - P(y=1 | x; \theta)$$

Multiclass Classification —

Likelihood function

$$\{x^{(i)}, y^{(i)}\} \quad ; \quad i = 1, 2, \dots, n$$

Each $y^{(i)}$ is IID.

cat
dog
horse

1
0
0

→ This corresponds to cat

$x_0^{(1)}$

$x_1^{(1)}$

$x_2^{(1)}$

$\rightarrow y^{(1)}$

1
0
0

{cat}

$x_0^{(2)}$

$x_1^{(2)}$

$x_2^{(2)}$

$\rightarrow y^{(2)}$

0
1
0

{dog}

$$f(y^{(i)}; \phi) = \phi_1^{\{y^{(i)}=1\}} \cdot \phi_2^{\{y^{(i)}=2\}} \cdots \phi_k^{\{y^{(i)}=k\}}$$

$$= \prod_{l=1}^k \phi_l^{\{y^{(i)}=l\}}$$

} Multinomial Distribution

Joint Probability

$$\hookrightarrow f(y, \phi) = f(y^{(1)}, \phi) \cdot f(y^{(2)}, \phi) \cdots \cdot f(y^{(n)}, \phi)$$

$$= \prod_{i=1}^n f(y^{(i)}, \phi)$$

$$= \prod_{i=1}^n \prod_{l=1}^k \phi_l^{\{y^{(i)}=l\}}$$

$$\phi_j = P(y^{(i)}=j) = P(y^{(i)}=j | x_i; \theta) = \frac{e^{\theta_j^T x^{(i)}}}{\sum_{s=1}^k e^{\theta_s^T x^{(i)}}}$$

Subject:

$$f(y, \theta) = \prod_{i=1}^n \prod_{l=1}^k \left(\frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)^{1\{y^{(i)}=l\}}$$

$L(\theta) \rightarrow$ likelihood function

log likelihood function = $\ell(\theta)$

$\log(L(\theta)) = \ell(\theta)$

$$= \log \left(\prod_{i=1}^n \prod_{l=1}^k \left(\frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)^{1\{y^{(i)}=l\}} \right)$$

$$= \sum_{i=1}^n \log \prod_{l=1}^k \left(\frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)^{1\{y^{(i)}=l\}}$$

$$= \sum_{i=1}^n \sum_{l=1}^k 1\{y^{(i)}=l\} \log \left(\frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)$$

Cost Function, $J(\theta) = -\ell(\theta)$

$$= - \sum_{i=1}^n \sum_{l=1}^k 1\{y^{(i)}=l\} \log P(y^{(i)}=l | x^{(i)}; \theta)$$



Cross entropy loss function
which is a convex function

$$h_{\theta}(x) = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \phi_3 \end{bmatrix} = \frac{1}{\sum_{j=1}^3 e^{\theta_j^T x}} \begin{bmatrix} \exp(\theta_1^T x) \\ \exp(\theta_2^T x) \\ \exp(\theta_3^T x) \end{bmatrix}$$

Till now we have considered the Learning Algorithms that modeled as $f(y|x; \theta)$ i.e., conditional distribution of the output y given the input feature vector x and parameterized by the weight vector θ .

$f(y|x; \theta) \sim \text{Normal}$ (Linear Regression)

$f(y|x; \theta) \sim \text{Bernoulli}$ (Logistic Regression)

$f(y|x; \theta) \sim \text{Multinomial}$ (Multiclass Regression)

$f(y|x; \theta) \sim \text{Exponential Distribution (GLMs)}$

These algorithms are called Discriminative Learning Algorithms.

We would now consider the algorithms that also consider/model the input feature vector x i.e., we would define $f(x|y; \theta)$. These algorithms are called Generative Learning Algorithms.

$f(x|y; \theta)$

↓
Conditional Probability density function
(pdf) of input feature vector x for
given y and θ .

Subject:

/ /

$P(y)$ → Prior probability of y

Using Bayes' Theorem;

$$P(y|x; \theta) = \frac{f(x|y; \theta) P(y)}{f(x)}$$

→ Discrete Random Variable

↓
Posterior distribution
of y given x and θ

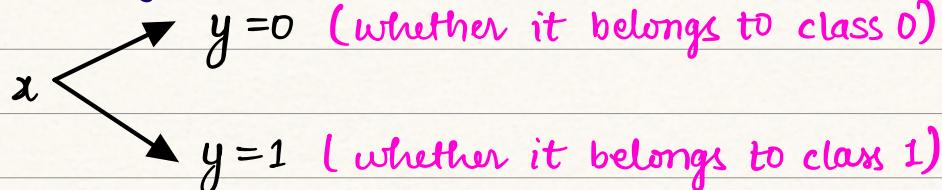
→ Continuous Random Variable

$$P(y|x; \theta) = \frac{f(x|y; \theta) P(y)}{f(x|y=0; \theta) P(y=0) + f(x|y=1; \theta) P(y=1)}$$

$\theta \rightarrow$ Known

Now if I give x as input then how will you decide either $y=0$ or $y=1$?

New training example $x = [x_0 \ x_1 \ x_2 \dots, x_N]$



To make decision, we should consider;

$$\arg \max_y P(y|x; \theta) = \arg \max_y f(x|y; \theta) P(y)$$

$y \in \{0, 1\}$

Generative Learning Algorithms :—

First is Gaussian Distribution Analysis (GDA)

Subject:

/ /

$f(\mathbf{x} | \mathbf{y}; \theta)$ is distributed according to a multivariate normal distribution.

↳ Input feature is vector.

* Multivariate Normal Distribution :—

$$\mathbf{x} \in \mathbb{R}^d$$

$$E[\mathbf{x}] = E \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} E[x_1] \\ E[x_2] \\ \vdots \\ E[x_d] \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{bmatrix} \in \mathbb{R}^d$$

(Mean vector, μ)

Covariance Matrix :

$$\text{Covariance } \Sigma = E \left[\underbrace{(\mathbf{x} - E[\mathbf{x}])}_{\mathbb{R}^{d \times 1}} \underbrace{(\mathbf{x} - E[\mathbf{x}])^T}_{\mathbb{R}^{1 \times d}} \right]$$

$$f(\mathbf{x}; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right)$$

For $d=1$;

$$\begin{aligned} \Sigma &= E[(\mathbf{x} - E[\mathbf{x}])(\mathbf{x} - E[\mathbf{x}])^T] \\ &= E[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] \\ &= E[(\mathbf{x} - \mu)(\mathbf{x} - \mu)] \\ &= E[(\mathbf{x} - \mu)^2] \\ &= \sigma^2 \end{aligned}$$

$$f(\mathbf{x}; \mu, \Sigma) = \frac{1}{(2\pi)^{1/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right)$$

Subject:

/ /

Suppose we have a classification problem in which the input features x are continuous valued and output y is discrete $\{0, 1\}$.

$$y \sim \text{Bernoulli}(\phi)$$
$$x | y=0 \sim \text{Multivariate Normal Distribution}$$
$$\sim N(\mu_0, \Sigma)$$

$$x = [x_0 \ x_1 \ x_2 \ \dots \ x_N]$$

$$P(y) = \phi^y (1-\phi)^{1-y}$$

$$P(y=1) = \phi \quad y \in \{0, 1\}$$
$$P(y=0) = 1 - \phi$$

$$f(x | y=0) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right)$$

$$f(x | y=1) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right)$$

Decision Rule;

$$\arg \max_y f(x | y; \theta) P(y)$$

$$\begin{cases} f(x | y=0; \theta) P(y=0) \\ f(x | y=1; \theta) P(y=1) \end{cases}$$

Parameters in Multivariate Normal distribution

$$f(x | y=0; \mu_0, \Sigma) P(y=0)$$
$$f(x | y=1; \mu_1, \Sigma) P(y=1)$$

Subject:

/ /

$$\begin{aligned} P(y=0) &= 1-\phi \\ P(y=1) &= \phi \end{aligned}$$

$$\mu_0 \text{MLE} = \frac{\sum_{i=1}^n 1\{y^{(i)}=0\} \cdot x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)}=0\}}$$

Maximum Likelihood
Estimator of μ_0

$$\mu_1 \text{MLE} = \frac{\sum_{i=1}^n 1\{y^{(i)}=1\} \cdot x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)}=1\}}$$

$$\phi = P(y=1)$$

$$\phi_{MLE} = \frac{\sum_{i=1}^n 1\{y^{(i)}=1\}}{n} = P(y=1)$$

Covariance Matrix $\Sigma = E[(x-\mu)(x-\mu)^T]$

$$= \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

$$\mu_1 \text{MLE} = \frac{\sum_{i=1}^n 1\{y^{(i)}=1\} \cdot x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)}=1\}} \quad \text{if } y^{(i)}=1$$

$$\mu_0 \text{MLE} = \frac{\sum_{i=1}^n 1\{y^{(i)}=0\} \cdot x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)}=0\}} \quad \text{if } y^{(i)}=0$$

Q.) which would perform better out of Gaussian Distribution Analysis (GDA) and logistic Regression?

Subject:

22/09/2025

$$\Sigma_{MLE} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}}) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T$$

Relation between Gaussian Discriminant Analysis (GDA)
and logistic Regression

$$f(\mathbf{x}|y) \sim \mathcal{N}(\text{Multivariate Gaussian})$$

$$P(y=1|\mathbf{x}, \phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) = \frac{f(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma) P(y=1|\mathbf{x}, \phi)}{f(\mathbf{x})}$$

$$f(\mathbf{x}) = P(y=1|\mathbf{x}; \phi) f(\mathbf{x}, \boldsymbol{\mu}_1, \Sigma) + P(y=0|\mathbf{x}; \phi) f(\mathbf{x}, \boldsymbol{\mu}_0, \Sigma)$$

$$= \frac{P(y=1|\mathbf{x}; \phi) f(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma)}{P(y=1|\mathbf{x}; \phi) f(\mathbf{x}, \boldsymbol{\mu}_1, \Sigma) + P(y=0|\mathbf{x}; \phi) f(\mathbf{x}, \boldsymbol{\mu}_0, \Sigma)}$$

$$= \frac{\phi \mathcal{N}(\mathbf{x}, \Sigma)}{\phi \mathcal{N}(\mathbf{x}, \Sigma) + (1-\phi) \mathcal{N}(\mathbf{x}, \Sigma)}$$

$$= \frac{1}{1 + (1-\phi) \frac{\mathcal{N}(\mathbf{x}, \Sigma)}{\phi \mathcal{N}(\mathbf{x}, \Sigma)}}$$

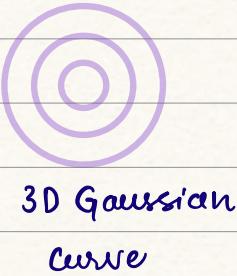
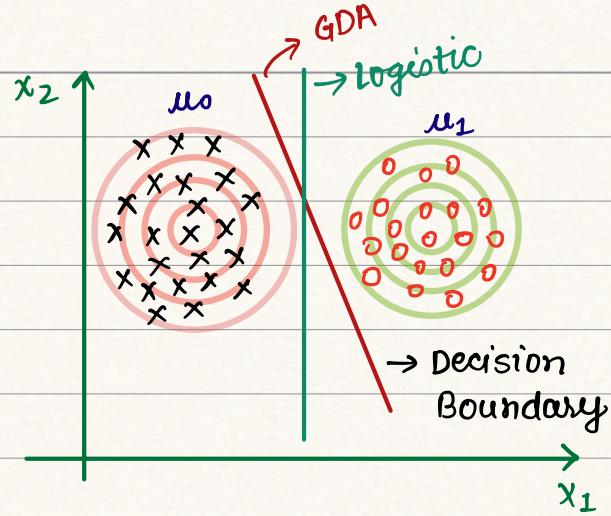
$$\mathcal{N}(\boldsymbol{\mu}_0, \Sigma) = \frac{1}{2\pi^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_0)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_0)\right)$$

$$\mathcal{N}(\boldsymbol{\mu}_1, \Sigma) = \frac{1}{2\pi^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)\right)$$

When we assume \mathbf{x} to be multivariate Gaussian, then

$$P(y=1|\mathbf{x}, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \phi, \Sigma) = \text{sigmoid}$$

Subject:



$$\begin{array}{ll} x_1^{(1)} & x_2^{(2)} \\ x_1^{(2)} & x_2^{(2)} \\ x_1^{(3)} & x_2^{(3)} \end{array} \quad \begin{array}{ll} y = 0 \\ y = 1 \\ y = 0 \end{array}$$

Q1) When would we prefer GDA or logistic regression?

$f(x|y) \sim \text{Multivariate Gaussian}$

$P(y=1|x) = \text{sigmoid function}$

But converse is not true.

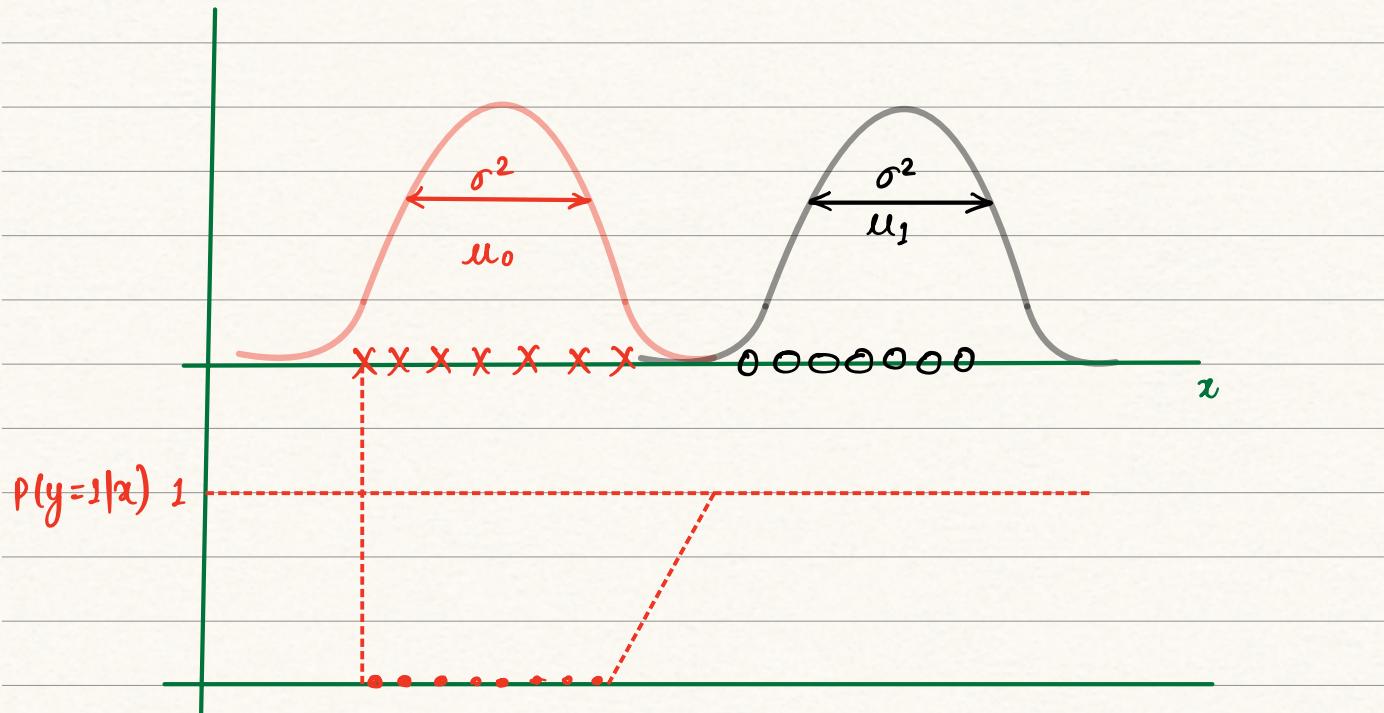
$P(y=1|x) \rightarrow \text{sigmoid function does not necessarily imply } P(x|y) \text{ is multivariate gaussian.}$

\Rightarrow GDA has stronger modelling assumptions as compared to logistic regression.

Data is
→ Gaussian

If the Gaussian modelling assumptions are correct, then GDA will find better performance as compared to logistic regression. Else GDA would perform poor.

$P(y=1|x) = \text{Sigmoid}$



* Kernel Method —

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

$$\phi: \mathbb{R} \rightarrow \mathbb{R}^4$$

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$y = \theta^T \phi(x)$$

↳ Complexity increased

that's why we are using kernel

In the context of kernels, the original input ' x ' is called Input attribute.

$\phi(x) \rightarrow$ Feature map which maps the input attribute (x) to features.

Subject:

/ /

update rule with $\phi(x)$;

$$h(x^{(i)}) = \theta^T \phi(x^{(i)})$$

$$\theta_j := \theta_j + \frac{\alpha}{n} \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) \cdot \phi(x_j^{(i)})$$

Generalized;

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$$

ϕ takes the input vector x in \mathbb{R}^d .

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_{d-1} \end{bmatrix} \quad (\mathbb{R}^d) \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{p-1} \end{bmatrix} \quad (\mathbb{R}^p)$$

Suppose that $d = 10^3$ then $p = d^3 = 10^9$. Since 10^9 is huge, as it will take lot of time to update 10^9 parameters.

Instead of 10^3 , the algorithm is now more complex.

* Least Mean Squared (LMS) update Rule with kernels :—

Q) How will we update θ with kernels?

Our claim: At any time, θ can be expressed as a linear combination of $\phi(x^{(1)})$, $\phi(x^{(2)})$, ... $\phi(x^{(n)})$

$$\theta = a_1 \phi(x^{(1)}) + a_2 \phi(x^{(2)}) + \dots + a_n \phi(x^{(n)})$$

$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

Suppose we initially set $\theta = 0$ (zero vector).

$$\theta := 0 + \alpha \sum_{i=1}^n (y^{(i)} - 0) \cdot \phi(x^{(i)})$$

New updated weight, $\theta = \sum_{i=1}^n \alpha \cdot y^{(i)} \cdot \phi(x^{(i)})$ } 2nd iteration

$$\theta = \sum_{i=1}^n \beta^{(i)} \cdot \phi(x^{(i)})$$

where $\beta^{(i)} = \alpha \cdot y^{(i)}$

$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

$$\theta = \sum_{i=1}^n \beta^{(i)} \cdot \phi(x^{(i)}) + \alpha \sum_{i=1}^n \left[y^{(i)} - \left(\sum_{j=1}^n \beta_j \cdot \phi(x^{(j)}) \right)^T \cdot \phi(x^{(i)}) \right] \cdot \phi(x^{(i)})$$

$$= \sum_{i=1}^n \beta_i \cdot \phi(x^{(i)}) + \alpha \sum_{i=1}^n \left(y^{(i)} - \sum_{j=1}^n \beta_j \cdot \phi(x^{(j)})^T \cdot \phi(x^{(i)}) \right) \cdot \phi(x^{(i)})$$

$$= \sum_{i=1}^n \left(\beta_i + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j \cdot \phi(x^{(j)})^T \cdot \phi(x^{(i)}) \right) \right) \cdot \phi(x^{(i)})$$

$$\beta_i^{(\text{new})} := \beta_i + \alpha \left(y^{(i)} - \underbrace{\sum_{j=1}^n \beta_j \cdot \phi(x^{(j)})^T \cdot \phi(x^{(i)})}_{O(np)} \right)$$

Complexity does not improved

Example :

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

$\phi(x)$ which contains the monomials of x with degree ≤ 3 .

Subject:

$$\underbrace{\langle \phi(x), \phi(z) \rangle}_{= 1 + \sum_{i=1}^d x_i z_i + \sum_{i,j} x_i x_j z_i z_j + \sum_{i,j,k} x_i x_j x_k z_i z_j z_k}$$

/ /

Dot product

$i, j \in [1, \dots, d]$

$i, j, k \in [1, \dots, d]$

$$= 1 + \underbrace{\langle x, z \rangle}_{= K(x, z)} + (\langle x, z \rangle)^2 + (\langle x, z \rangle)^3$$

$$= K(x, z) \quad O(d)$$

$d \ll p$

This is Kernel.

$$\langle \phi(x), \phi(z) \rangle = K(x, z) = \text{kernel}$$

$$\beta_i^{\text{new}} := \beta_i + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j \underbrace{\phi(x^{(j)})^T \phi(x^{(i)})}_{\text{This is fixed for each training example.}} \right)$$

$$= \beta_i + \alpha \left(y^{(i)} - \sum_{j=1}^n \beta_j \underbrace{K(x^{(j)}, x^{(i)})}_{O(d)} \right)$$

$O(nd)$

$$K(x^{(j)}, x^{(i)}) = \langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle$$

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

As this is fixed for each training example

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

we will store its value separately so

$$i, j = 1, 2, \dots, n$$

that algorithm becomes much more

faster.

Let suppose

Q) What kind of function $K(\cdot)$ can corresponds to some feature map, ϕ ? In other words, can we tell if there is some feature mapping ϕ so that $K(x, z) = \phi(x)^T \phi(z)$

Example :-

$$x, z \in \mathbb{R}^d$$

$$K(x, z) = (x^T z)^2 \rightarrow \text{polynomial kernel of degree 2}$$

$$x = [x_1 \ x_2 \ \dots \ x_d]$$

$$z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_d \end{bmatrix}$$

$$\begin{aligned} K(x, z) &= (x^T z)^2 \\ &= (x^T z)(x^T z) \\ &= \left(\sum_{i=1}^d x_i z_i \right) \left(\sum_{j=1}^d x_j z_j \right) \end{aligned}$$

$$\begin{aligned} &= \sum_{i=1}^d \sum_{j=1}^d x_i x_j z_i z_j \\ &= \langle \phi(x), \phi(z) \rangle \end{aligned}$$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_2 x_1 \\ x_2 x_2 \end{bmatrix}$$

$$\phi(z) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_2 z_1 \\ z_2 z_2 \end{bmatrix}$$

Subject:

/ /

Suppose $d = 3$;

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \quad \phi(z) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix}$$

$$K(x, z) = (x^T z + c)^2$$

↳ Polynomial kernel of degree 2

$$= \sum_{i,j=1}^d (x_i x_j) \cdot (z_i z_j) + \sum_{i=1}^d (\sqrt{c} x_i) (\sqrt{c} z_i) + c^2$$

$$= \langle \phi(x), \phi(z) \rangle$$

$d=3$

$$\phi(x) =$$

$$\begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ \vdots \\ x_3 x_3 \\ \sqrt{c} \cdot x_1 \\ \sqrt{c} \cdot x_2 \\ \sqrt{c} \cdot x_3 \\ c \end{bmatrix}$$

$$\phi(z) =$$

$$\begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ \vdots \\ z_3 z_3 \\ \sqrt{c} \cdot z_1 \\ \sqrt{c} \cdot z_2 \\ \sqrt{c} \cdot z_3 \\ c \end{bmatrix}$$

★ Necessary condition for a valid kernel —

Actual inputs = $[x^{(1)}, x^{(2)}, \dots x^{(n)}]$

$$K_{ij} = K(x^{(i)}, x^{(j)})$$

$i, j \in [1, 2, \dots n]$

$$K = \text{Kernel matrix} = \begin{bmatrix} K_{11} & K_{12} & \dots & K_{1n} \\ K_{21} & K_{22} & \dots & K_{2n} \\ \vdots & \vdots & \dots & \vdots \\ K_{n1} & K_{n2} & \dots & K_{nn} \end{bmatrix}_{n \times n}$$

For a Kernel $K(\cdot)$ to be valid,

1. $K_{ij} = K_{ji} \rightarrow$ Symmetric matrix

2. For any $z \in \mathbb{R}^n$,

$$z^T K z \geq 0$$

$K \rightarrow$ Positive semi Definite

★ Gaussian Kernel or Radial basis function (RBF) Kernel :-

$x, z \in \mathbb{R}^d$

$$K(x, z) = \langle \phi(x), \phi(z) \rangle$$



Very large when
 x and z are very
close to each other

Very small when
 x and z are far apart
or orthogonal to each other

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2}\right) = \langle \phi(x), \phi(z) \rangle$$

Subject:

/ /

* It has infinite dimensions feature map.

Proof: Gaussian Kernel : $K(u, v) = \exp\left(-\frac{\|u-v\|^2}{2}\right)$

$$\begin{aligned}\|u-v\|^2 &= (u-v)^T(u-v) \\&= (u^T - v^T)(u-v) \\&= u^Tu - u^Tv - v^Tu + v^Tv \\&= \|u\|^2 - 2u^Tv + \|v\|^2\end{aligned}$$

$$\exp\left(-\frac{\|u-v\|^2}{2}\right) = \exp\left(-\frac{1}{2}(\|u\|^2 - 2u^Tv + \|v\|^2)\right)$$

$$= \exp\left(-\frac{\|u\|^2}{2}\right) \exp\left(-\frac{\|v\|^2}{2}\right) \exp(\langle u, v \rangle)$$

$$\exp(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

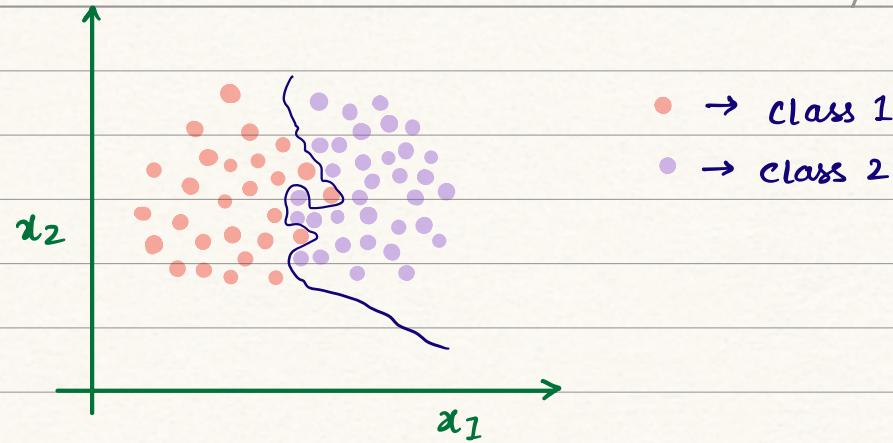
$$= \exp\left(-\frac{\|u\|^2}{2}\right) \exp\left(-\frac{\|v\|^2}{2}\right) \sum_{k=0}^{\infty} \left(\frac{\langle u, v \rangle}{k!}\right)^k$$

* Neural Network —

Q) Why do we need Neural Network when we already have the algorithms such as Linear Regression and Logistic Regression?

Suppose we have a training set corresponding to Logistic Regression when we need to learn complex non-linear decision-boundary.

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

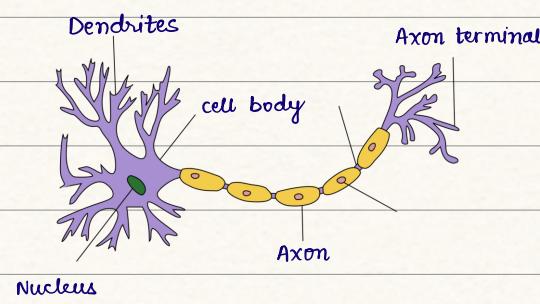


In real-life, there are more than 30 input features and to train on the dataset, the complexity will increase exponentially. So that is why we use neural networks to solve this problem.

- * Simple Linear/Logistic Regression together with a lot of quadratic or cubic features is not a good way to learn the complex non-linear decision boundary as it will lead to huge complexity.
- * We should consider the neural networks which turn out to be a much better way to learn complex non-linear decision boundaries with less complexity even when the number of inputs are very large.

Neural Networks :- Algorithm that try to mimic the human brain

Biological Neuron —

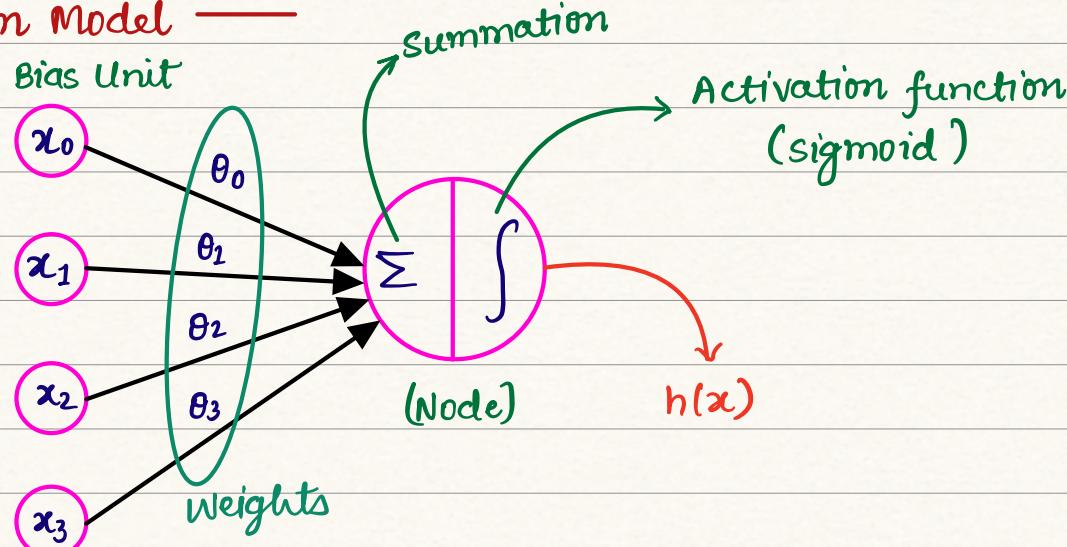


Subject:

/ /

Biological Neuron	Artificial Neuron
Dendrites	Inputs
Cell Nucleus	Node
Synapsis	Weights
Axon	Outputs

Neuron Model —

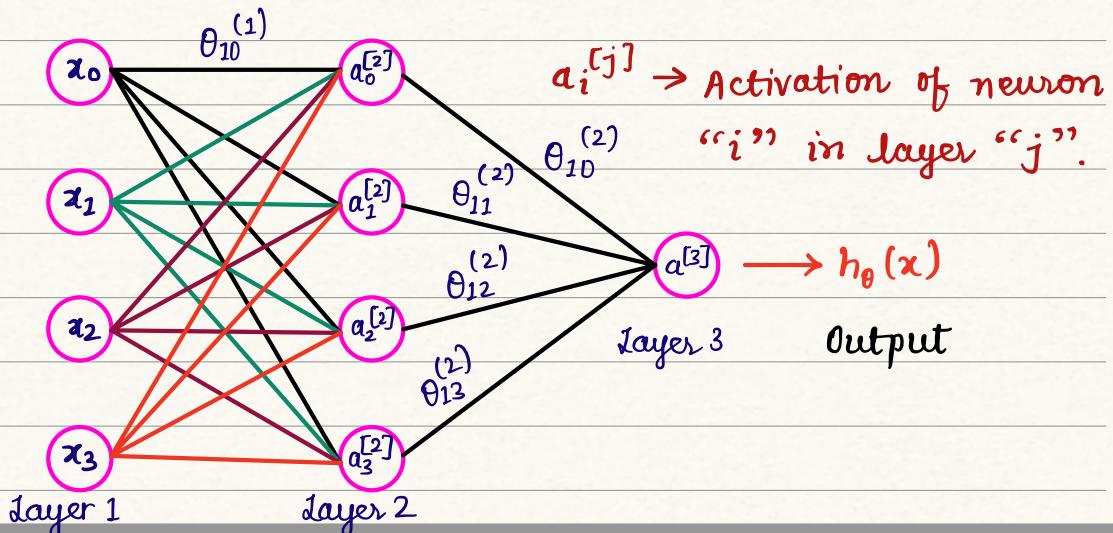


$$\sum = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$\int = \frac{1}{e^{-\theta^T x}} = \frac{1}{e^{-(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3)}}$$

Neural Network —

It is the group of the different neurons stuck together.



Subject:

Feed

Forward

Propagation

layer 1 :- Input Layer

Layer 2 :- Hidden Layer

Layer 3 :- Output Layer

Backward

Propagation

$\theta^{(j)}$

Matrix/Vector of weights controlling function mapping from layer ' j ' to layer ' $j+1$ '.

$$a_i^{[j]} = \sum \int$$

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

$g \rightarrow$ sigmoid activation function

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

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

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

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

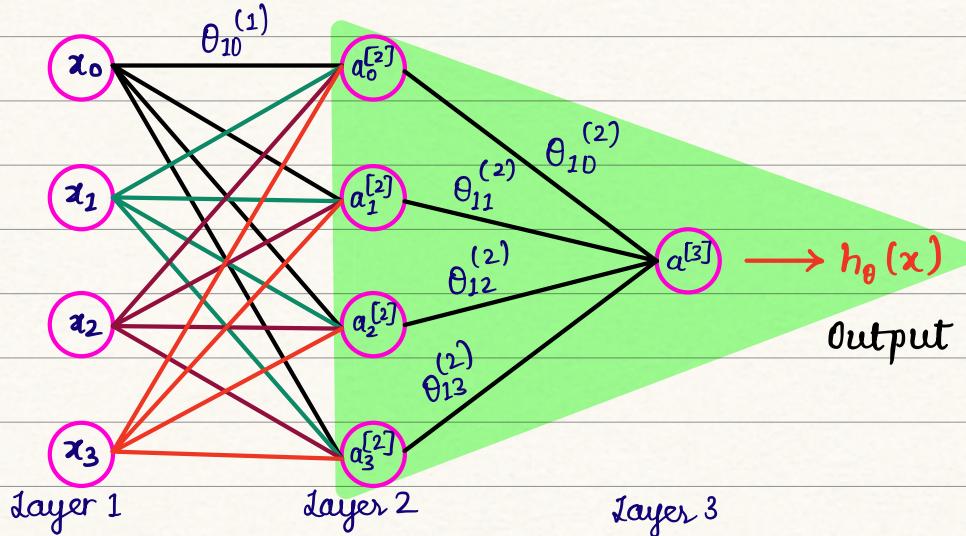
$$\theta^{(1)} = \begin{bmatrix} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \\ \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \\ \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{bmatrix}$$

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

$$\theta^{(1)} x = z^{[2]} \quad \theta^{(1)} \in \mathbb{R}^{3 \times 4}$$
$$(3 \times 4) \quad (4 \times 1) \quad (3 \times 1)$$

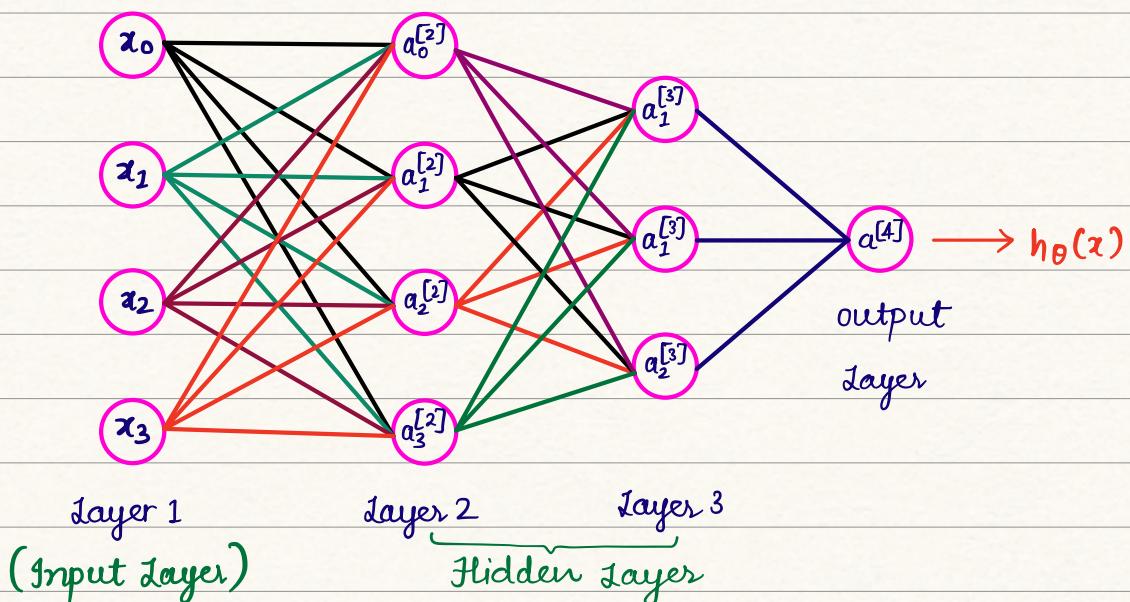
$$a^{[3]} = g(a_0^{[2]} \theta_{10}^{(2)} + a_1^{[2]} \theta_{11}^{(2)} + a_2^{[2]} \theta_{12}^{(2)} + a_3^{[2]} \theta_{13}^{(2)})$$

* Neural Network learning its own features —

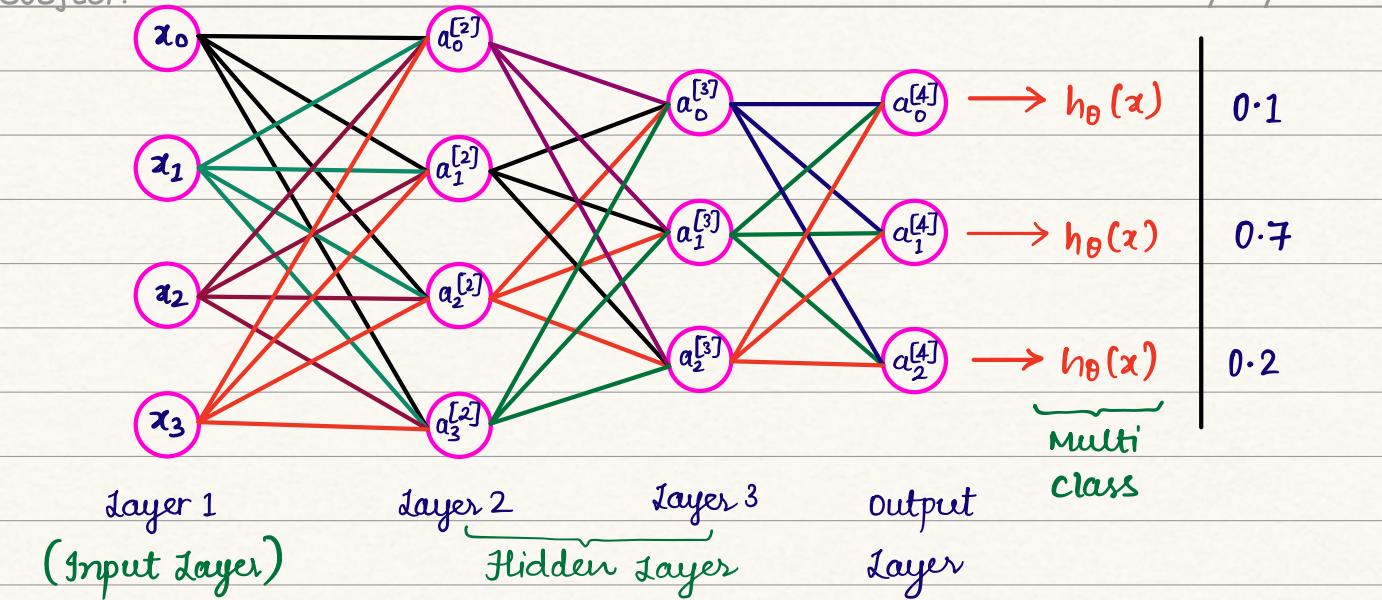


$$h_\theta(x) = g(\theta_{10}^{(2)} a_0^{[2]} + \theta_{11}^{(2)} a_1^{[2]} + \theta_{12}^{(2)} a_2^{[2]} + \theta_{13}^{(2)} a_3^{[2]})$$

* In this Neural network Architecture, it gets to learn its own features $a_1^{[2]}, a_2^{[2]}, a_3^{[2]}$ to feed into the logistic regression. With this, with the help of Neural Networks, we can easily learn complex non-linear decision boundaries.



Subject:



Output Layer : $a^{[j]} \rightarrow$ Activation function

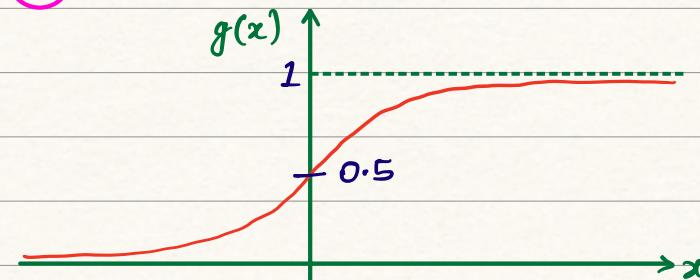
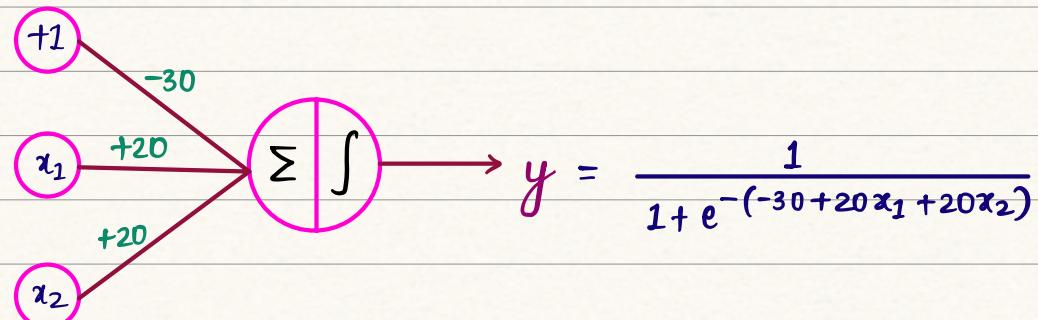
$$\text{House Price} \\ h_\theta(x) = \theta^T x$$

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

* Logic Gates using Neural Networks —

$x_1, x_2 \in \{0, 1\} \rightarrow$ Binary Inputs

$y = x_1 x_2 = \text{AND operation of } x_1 \text{ and } x_2$



Subject:

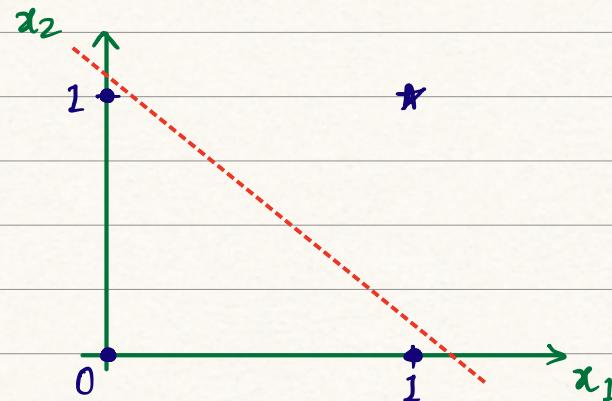
/ /

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

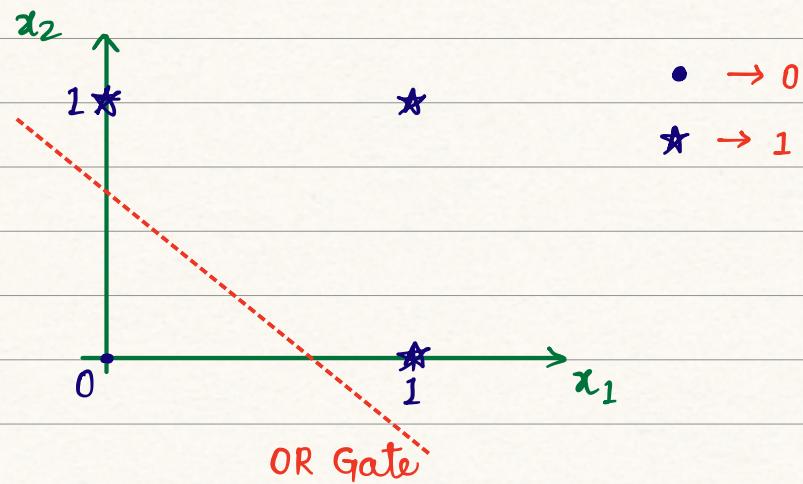
AND Gate

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1

OR Gate

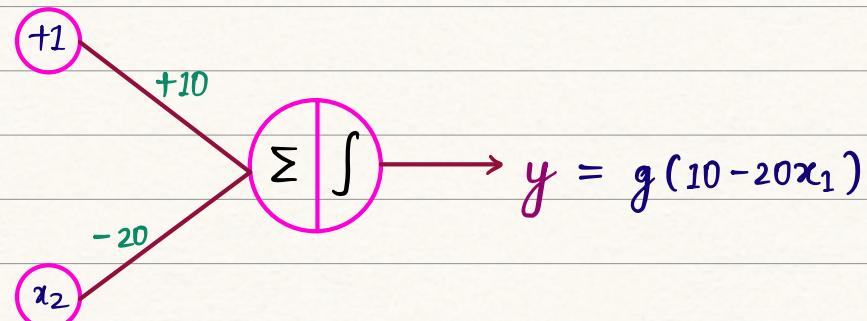


AND Gate



OR Gate

NOT Gate :

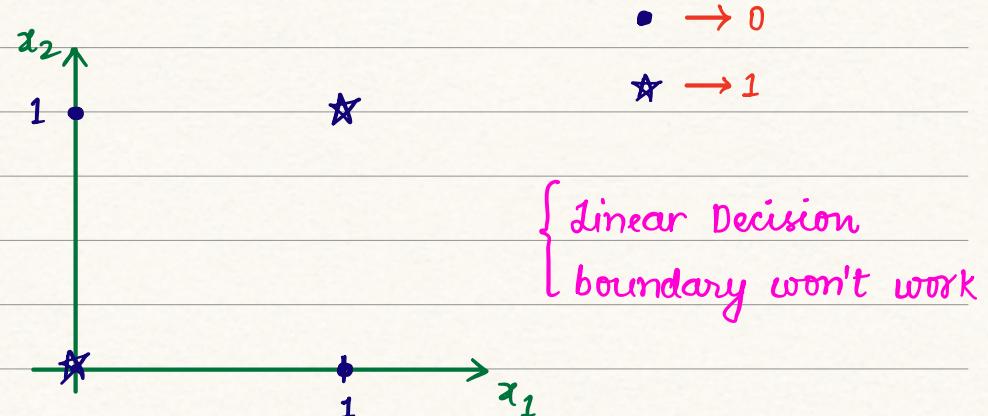


Subject:

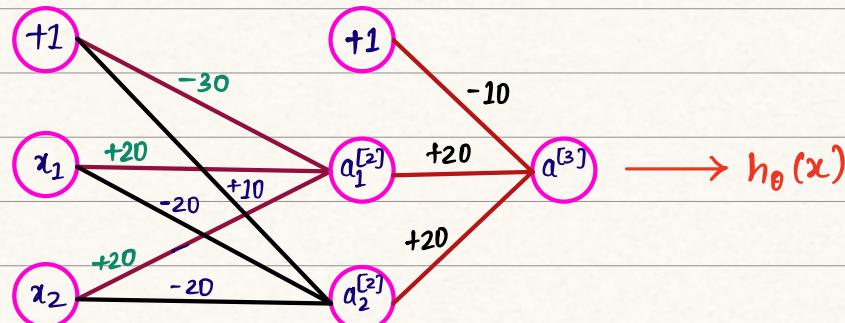
/ /

* Non-linear classification Example —

$$XNOR = x_1 x_2 + \bar{x}_1 \bar{x}_2$$



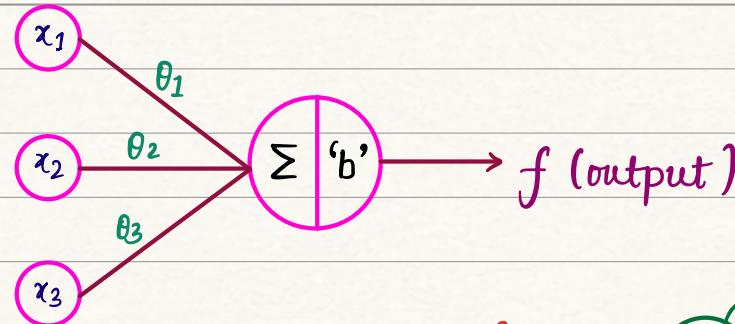
$$x_1 x_2 + \bar{x}_1 \bar{x}_2$$



x_1	x_2	$a_1^{[2]}$	$a_2^{[2]}$	$h_\theta(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

* Perceptron —

Take several inputs and gives a single output (binary).



$$f = \begin{cases} 0 & \text{if } \sum_{i=1}^3 x_i \theta_i \leq -b \\ 1 & \text{if } \sum_{i=1}^3 x_i \theta_i > -b \end{cases}$$

Threshold

$$f = \begin{cases} 0 & \text{if } \theta^T x + b \leq 0 \\ 1 & \text{if } \theta^T x + b > 0 \end{cases}$$

$$\text{def } z = \theta^T x + b$$

$$\text{Activation, } a = f(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases}$$



* Perceptron Example —

Q.) Will you go to the picnic ?

x_1 : Is the weather good ?

x_2 : Do you have any company with friends ?

x_3 : Is there any good garden nearby ?

Subject:

/ /

We decide that we will go to picnic when the weather is good and if any one of the rest of the factors is in my favour.

"Weather"

"Good company"

"Good Garden nearby"

$$x_1 \quad \theta_1 = 4$$

$$x_2 \quad \theta_2 = 2$$

$$\theta_3 = 2$$

$$\Sigma b \rightarrow f$$

$$\{ b = -5 \}$$

$$f = 4x_1 + 2x_2 + 2x_3 - 5$$

Suppose that $x_1 = 0$ (weather is bad)

$$x_1 = 0, x_2 = 1, x_3 = 1$$

$$f = -1 < 0 \rightarrow 0$$

★ Perception also work on linear classifier —

$$x_1 \quad -2$$

$$x_2 \quad -2$$

$$\Sigma 'b' \rightarrow f$$

$$\{ b = 3 \}$$

$$f = -2x_1 - 2x_2 + 3$$

x_2

$3/2$

$f > 0$

$f \leq 0$

$3/2$

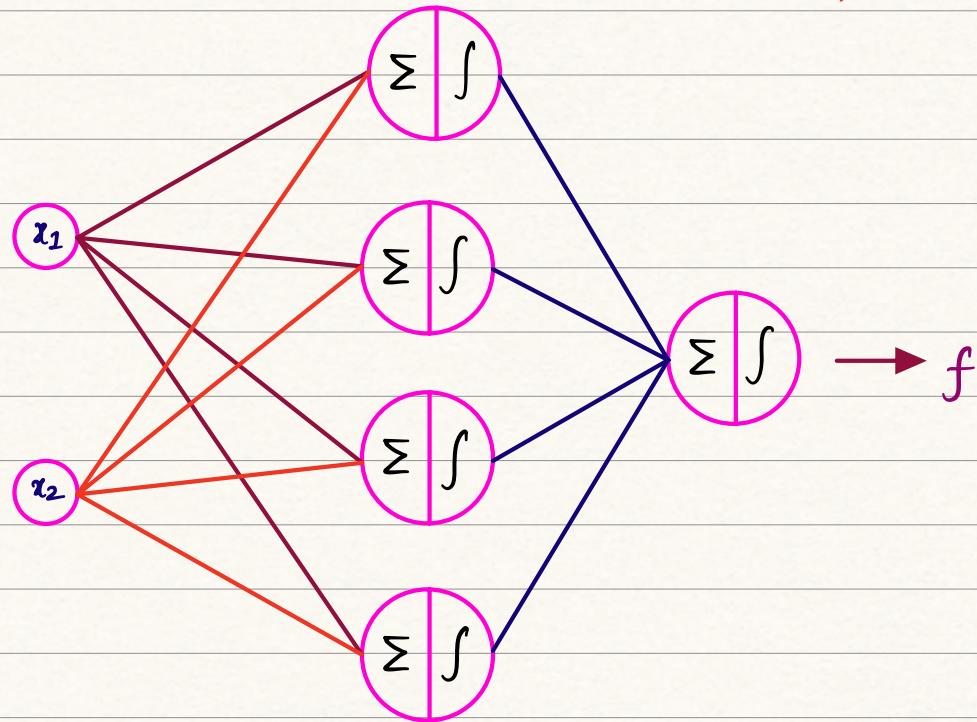
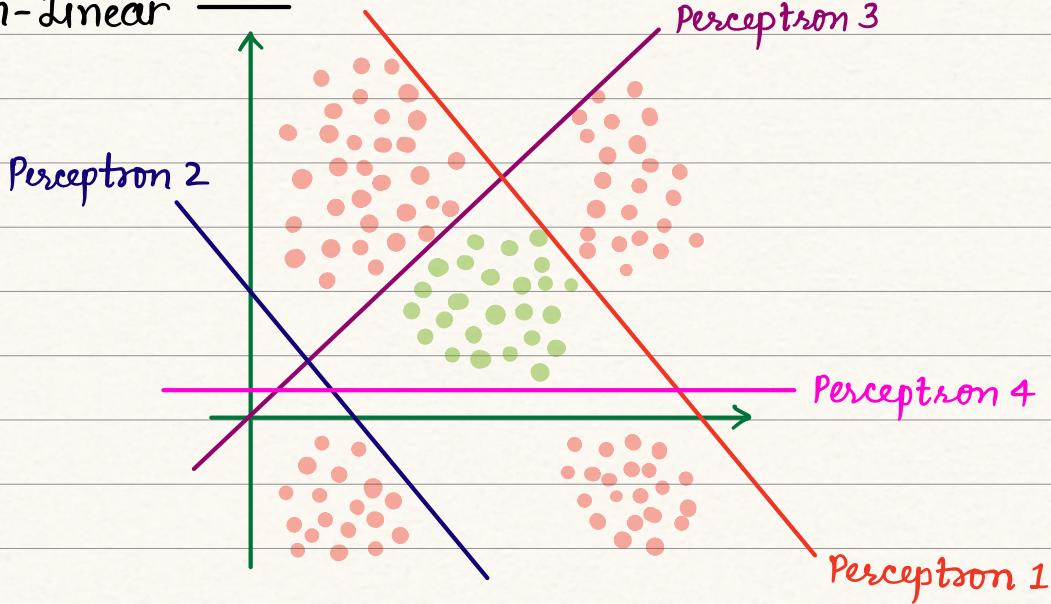
x_1

Subject: / /

x_1	x_2	f	$a(f)$
0	0	3	1
0	1	1	1
1	0	1	1
1	1	-1	0

NAND Gate

* Non-Linear



* Drawbacks of perception —

- (i) Step function
- (ii) Not differentiable

Subject:

/ /

Housing Price Prediction Problem :—

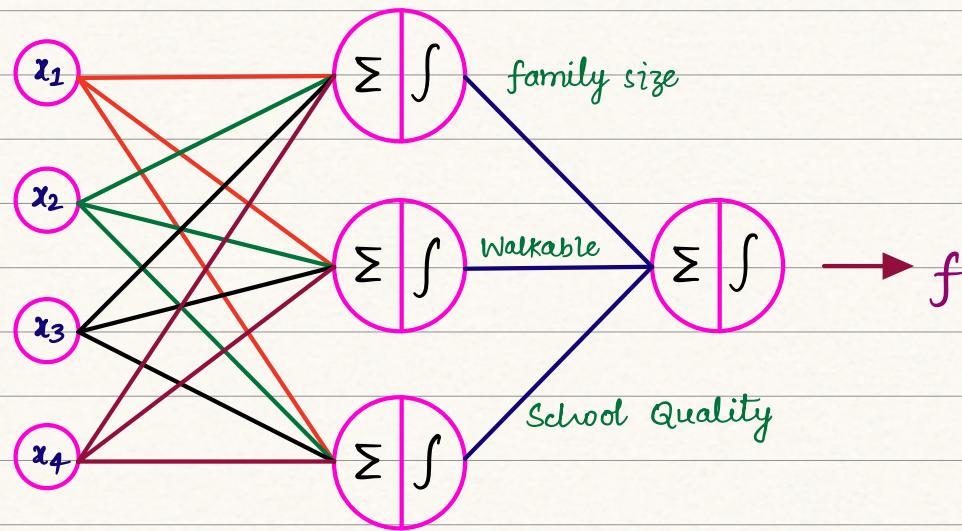
Suppose the inputs :

x_1 : size of house

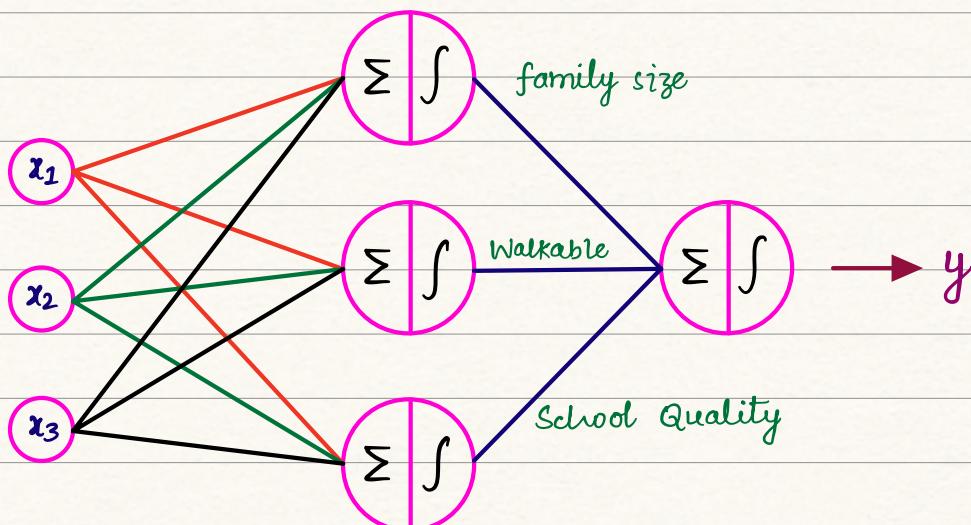
x_2 : number of bedrooms

x_3 : zip code

x_4 : wealth



Q.) Why we connect each input to each neuron ?



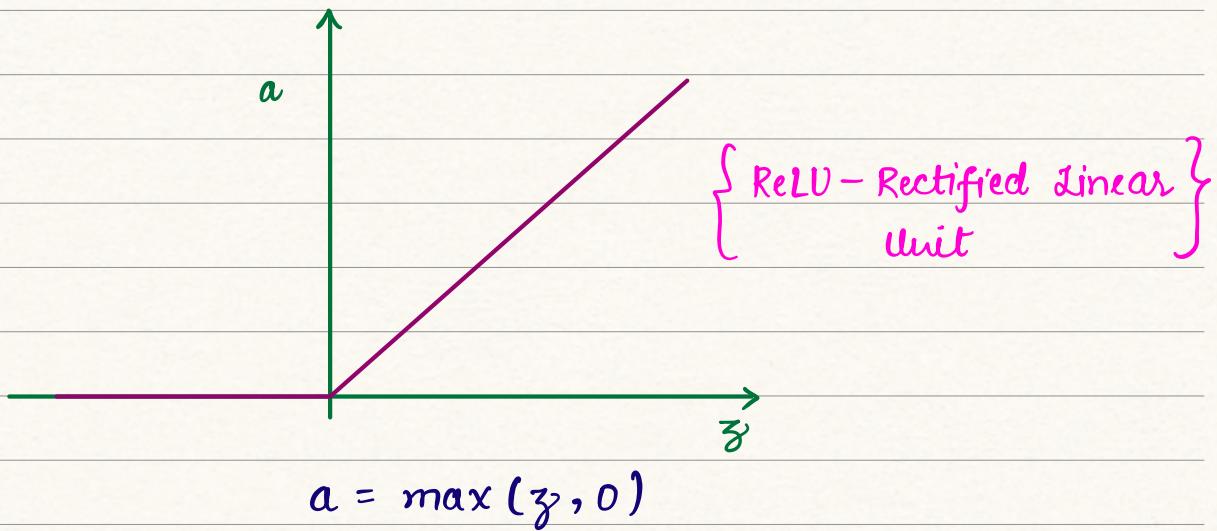
Subject:

/ /

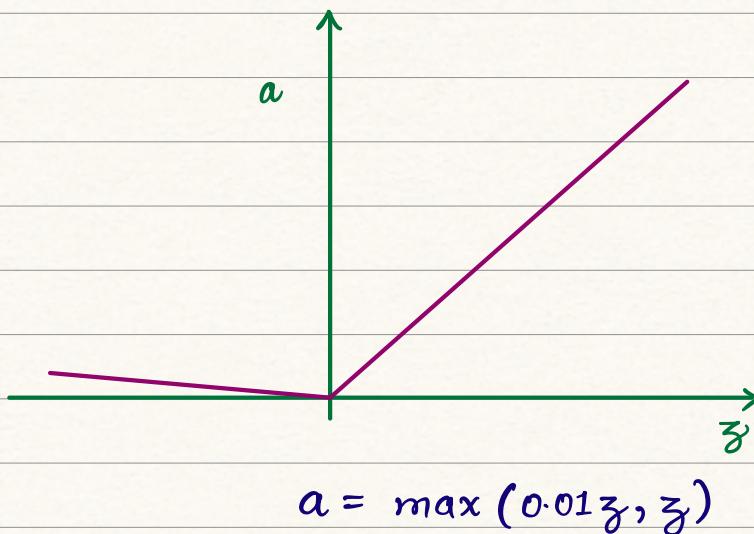
$$-1 \leq y \leq 1 \rightarrow \text{tanh Activation Function}$$

Q.) What is problem associated with sigmoid function for higher value of 'z'?

Solution to this problem : ReLU Activation Function

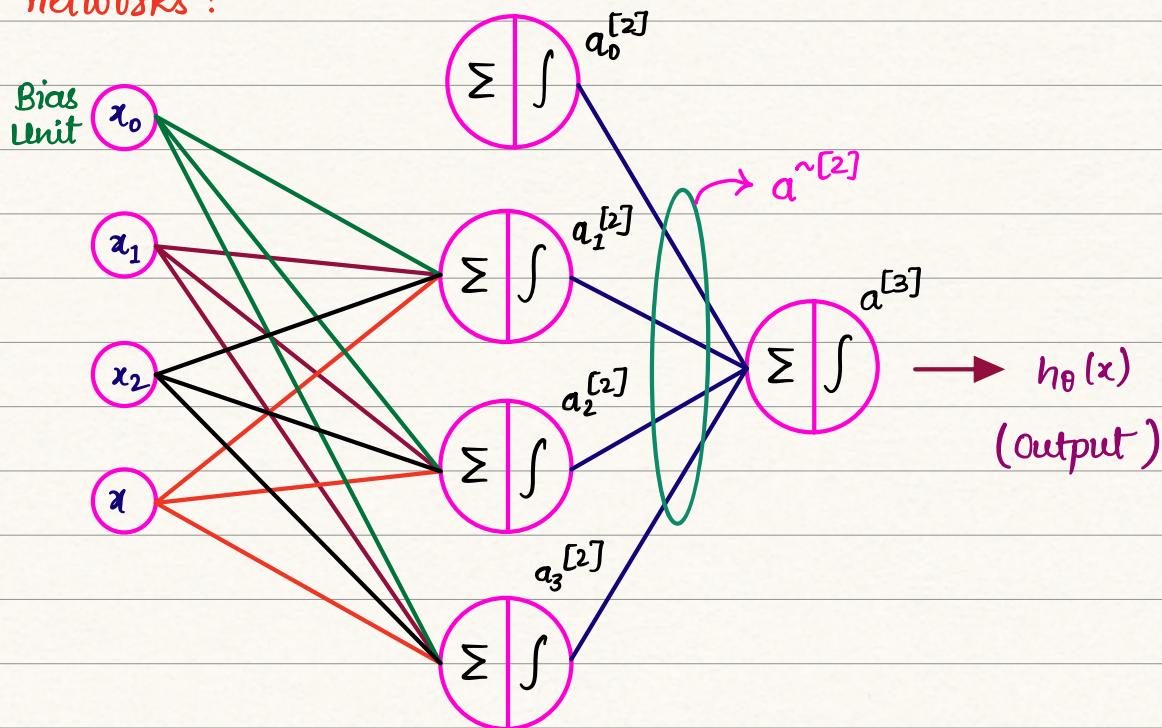


Leaky ReLU :—



$g \rightarrow$ Identity function
 $f(a) = a$

Q.) Why we need non-linear activation functions for neural networks?



Non-linear Activation Functions : —

- (i) Sigmoid
- (ii) tanh
- (iii) ReLU
- (iv) Leaky ReLU

$$\underline{z}^{[2]} = \underline{\theta}^{(1)} \underline{x} \in \mathbb{R}^{3 \times 1}$$

Matrix of weights while going from layer 1 to layer 2.

$$\underline{x} \in \mathbb{R}^{4 \times 1}$$

$$\underline{\theta} \in \mathbb{R}^{3 \times 4}$$

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

$$\underline{z}^{[3]} = \underline{\theta}^{(2)} \underline{a}^{[2]} \quad (1 \times 4) \quad (4 \times 1)$$

$$\underline{a}^{[2]} = \begin{bmatrix} a_0^{[2]} \\ a_1^{[2]} \\ a_2^{[2]} \\ a_3^{[2]} \end{bmatrix}_{4 \times 1}$$

$$\underline{a}^{[2]} = \begin{bmatrix} 1 \\ a_1^{[2]} \\ a_2^{[2]} \\ a_3^{[2]} \end{bmatrix} \rightarrow \underline{a}^{[2]}$$

Subject:

/ /

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

Instead of any Non-linear Activation Function;

$g(z) = z$ {Identity function}

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

$$z^{[2]} = \theta^{(1)} \underline{x}$$

$$a^{[2]} = g(z^{[2]}) = z^{[2]} = \theta^{(1)} \underline{x}$$

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

$$a^{[3]} = g(z^{[3]}) = z^{[3]} = \theta^{(2)} \underline{a}^{[2]}$$

$$= \theta^{(2)} \begin{bmatrix} a_0^{[2]} \\ \underline{a}^{[2]} \end{bmatrix}$$

$$\begin{cases} a_0^{[2]} = 1 \\ \underline{a}^{[2]} = z^{[2]} = \theta^{(1)} \underline{x} \end{cases}$$

$$= \theta^{(2)} \theta^{(1)} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\underline{a}^{[2]} = \begin{bmatrix} a_1^{[2]} \\ a_2^{[2]} \\ a_3^{[2]} \end{bmatrix}_{3 \times 1}$$

$$\theta^{(1)} = \begin{bmatrix} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \\ \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \\ \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{bmatrix}_{3 \times 4} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}_{4 \times 1}$$

$$\theta^{(2)} = \begin{bmatrix} \theta_{10}^{(2)} & \theta_{11}^{(2)} & \theta_{12}^{(2)} & \theta_{13}^{(2)} \end{bmatrix}$$

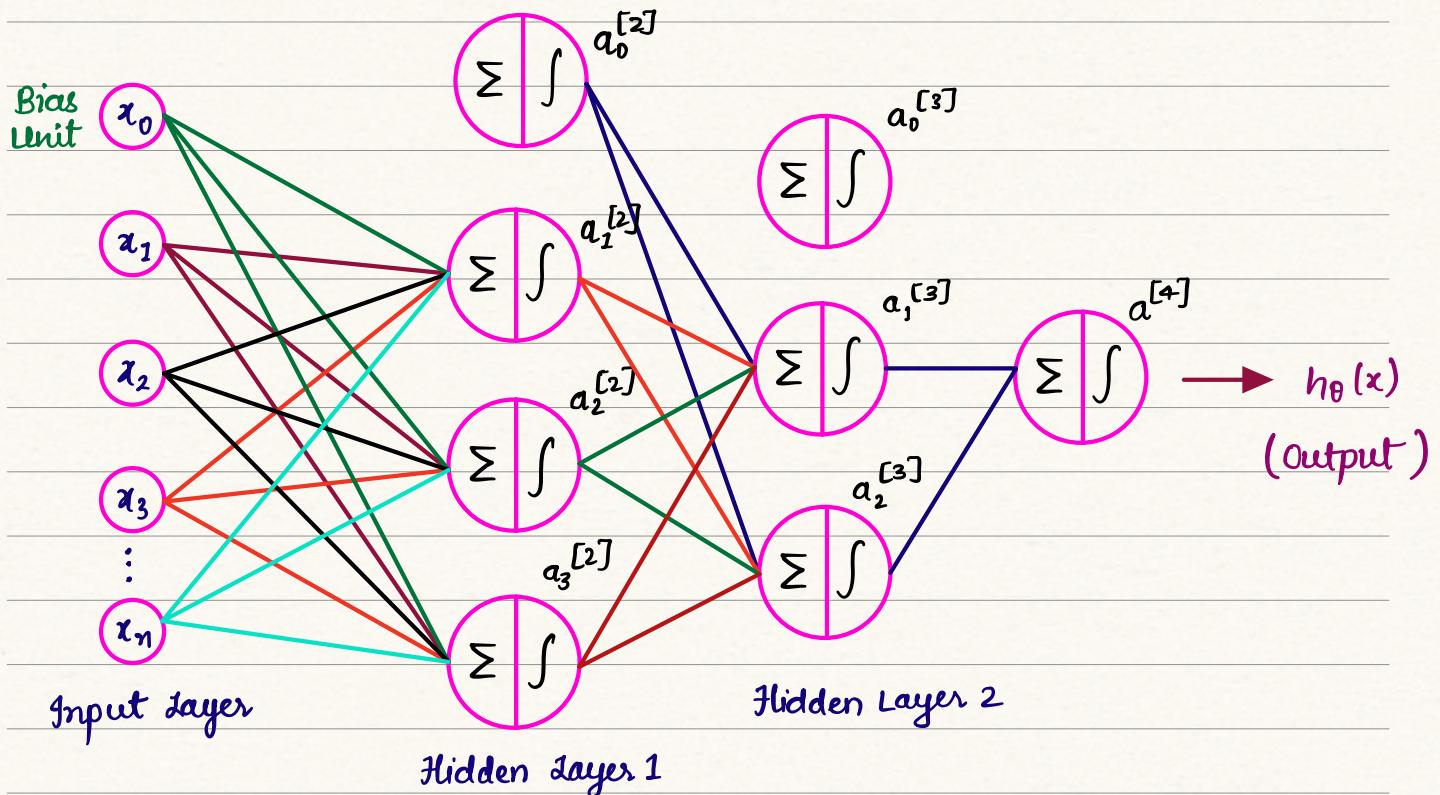
Subject:

$$z^{[2]} = \begin{bmatrix} \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3 \\ \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3 \\ \theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3 \end{bmatrix}$$

$$\begin{bmatrix} \theta_{10}^{(2)} & \theta_{11}^{(2)} & \theta_{12}^{(2)} & \theta_{13}^{(2)} \end{bmatrix} \begin{bmatrix} \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3 \\ \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3 \\ \theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3 \end{bmatrix} \rightarrow \begin{array}{l} a \\ b \\ c \end{array}$$

$$= \theta_{10}^{(2)} + \theta_{11}^{(2)} a + \theta_{12}^{(2)} b + \theta_{13}^{(2)} c$$

★ Backpropagation Algorithm —



$$z^{[1]} = \begin{bmatrix} \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3 \\ \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3 \\ \theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3 \end{bmatrix}$$

" $x_0 = 1$ ";

$$\vec{z}^{(1)} = \begin{bmatrix} \theta_{10}^{(1)} + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \theta_{13}^{(1)}x_3 \\ \theta_{20}^{(1)} + \theta_{21}^{(1)}x_1 + \theta_{22}^{(1)}x_2 + \theta_{23}^{(1)}x_3 \\ \theta_{30}^{(1)} + \theta_{31}^{(1)}x_1 + \theta_{32}^{(1)}x_2 + \theta_{33}^{(1)}x_3 \end{bmatrix}$$

$$= \begin{bmatrix} \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \\ \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \\ \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} \theta_{10}^{(1)} \\ \theta_{20}^{(1)} \\ \theta_{30}^{(1)} \end{bmatrix}$$

$$z^{(1)} = w^{(1)}x + b^{(1)}$$

$(3 \times 1) \quad (3 \times 4) \quad (4 \times 1)$

Bias Vector

$$a^{(1)} = \sigma(z^{(1)}) = \sigma(w^{(1)}x + b^{(1)})$$

$\{\sigma \rightarrow \text{sigmoid Activation function}\}$

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

$(2 \times 1) \quad (2 \times 3) \quad (3 \times 1)$

$$\begin{bmatrix} \theta_{10}^{(2)} & \theta_{11}^{(2)} & \theta_{12}^{(2)} & \theta_{13}^{(2)} \\ \theta_{20}^{(2)} & \theta_{21}^{(2)} & \theta_{22}^{(2)} & \theta_{23}^{(2)} \end{bmatrix} \begin{bmatrix} 1 \\ a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix}$$

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

$$z^{(3)} = w^{(3)}a^{(2)} + b^{(3)}$$

$(1 \times 1) \quad (1 \times 3) \quad (3 \times 1)$

$$a^{(3)} = \sigma(z^{(3)}) = \sigma(w^{(3)}a^{(2)} + b^{(3)})$$

Subject:

$$\begin{aligned} z^{(1)} &= w^{(1)}x + b^{(1)} \\ a^{(1)} &= \sigma(z^{(1)}) = \sigma(w^{(1)}x + b^{(1)}) \quad (3 \times 1) \end{aligned}$$

$$\begin{aligned} z^{(2)} &= w^{(2)}a^{(1)} + b^{(2)} \\ a^{(2)} &= \sigma(z^{(2)}) = \sigma(w^{(2)}a^{(1)} + b^{(2)}) \quad (2 \times 1) \end{aligned}$$

$$\begin{aligned} z^{(3)} &= w^{(3)}a^{(2)} + b^{(3)} \\ a^{(3)} &= \sigma(z^{(3)}) = \sigma(w^{(3)}a^{(2)} + b^{(3)}) \quad (1 \times 1) \end{aligned}$$

To train any Neural Network Model, we need to update the following parameters :

- | | |
|--------------|--------------|
| 1> $w^{(1)}$ | 4> $b^{(2)}$ |
| 2> $b^{(1)}$ | 5> $w^{(3)}$ |
| 3> $w^{(2)}$ | 6> $b^{(3)}$ |

Q.) How can I update these parameters ?

We need to define the loss function.

→ Optimize these parameters such that loss function is minimized.

Suppose that I have a dataset (training) that deals with the binary classification problems.

Suppose that I have "m" training examples.

Loss function = Binary cross Entropy loss function

$$= \sum_{i=1}^m - (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}))$$

For n = 1 ;

$$L = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

Subject:

/ /

$$w^{(l)} := w^{(l)} - \alpha \frac{\partial L}{\partial w^{(l)}} \quad l = \{1, 2, 3\}$$

$$b^{(l)} := b^{(l)} - \alpha \frac{\partial L}{\partial b^{(l)}}$$

$$L = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

$$\hat{y} = \sigma(z^{(3)}) = \sigma(w^{(3)}a^{(2)} + b^{(3)})$$

(1x2) (2x1)

$$\begin{aligned} \frac{\partial L}{\partial w^{(3)}} &= -y \frac{\partial}{\partial w^{(3)}} \log(\sigma(z^{(3)})) - (1-y) \frac{\partial}{\partial w^{(3)}} \log(1-\sigma(z^{(3)})) \\ &= \frac{-y}{\sigma(z^{(3)})} \sigma(z^{(3)})(1-\sigma(z^{(3)}))(a^{(2)})^T \\ &\quad - \frac{(1-y)}{(1-\sigma(z^{(3)}))} (1-\sigma(z^{(3)}))(-\sigma(z^{(3)}))(a^{(2)})^T \end{aligned}$$

$$x = \begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = w_1 a_1 + w_2 a_2$$

$$\frac{\partial x}{\partial w} = [a_1 \quad a_2]$$

$$= -(y(1-a^{(3)}) - (1-y)a^{(3)})(a^{(2)})^T$$

$$= -(y - a^{(3)}y - a^3 + ya^{(3)})(a^{(2)})^T \quad (a^{(2)})^T \in \mathbb{R}^{1 \times 2}$$

$$\frac{\partial L}{\partial w^{(3)}} = \underbrace{-(y - a^{(3)})}_{\frac{\partial L}{\partial z^{(3)}}} \underbrace{(a^{(2)})^T}_{\frac{\partial z^{(3)}}{\partial w^{(3)}}} \quad \left\{ (a^{(2)})^T = \frac{\partial z^{(3)}}{\partial w^{(3)}} \right.$$

$$L = -y \log(\sigma(z^{(3)})) - (1-y) \log(1-\sigma(z^{(3)}))$$

$$\begin{aligned}
 \frac{\partial L}{\partial z^{(3)}} &= \frac{-y}{\sigma(z^{(3)})} \sigma(z^{(3)}) (1 - \sigma(z^{(3)})) \\
 &\quad - \frac{(1-y)}{(1-\sigma(z^{(3)}))} (-\sigma(z^{(3)}) (1 - \sigma(z^{(3)}))) \\
 &= -y (1 - \sigma(z^{(3)})) + (1-y) (\sigma(z^{(3)})) \\
 &= -y + y \sigma(z^{(3)}) + \sigma(z^{(3)}) - y \sigma(z^{(3)}) \\
 &= -(y - \sigma(z^{(3)})) \\
 &= -(y - a^{(3)})
 \end{aligned}$$

$$\frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial z^{(3)}} \times \frac{\partial z^{(3)}}{\partial a^{(2)}} \times \frac{\partial a^{(2)}}{\partial z^{(2)}} \times \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

$\overbrace{(2 \times 3)}$ We are using chain rule because one layer has dependency on the previous layer.

$$L = -y \log a^{(3)} - (1-y) \log (1-a^{(3)})$$

$$\begin{aligned}
 a^{(3)} &= \sigma(z^{(3)}) \\
 z^{(3)} &= w^{(3)} a^{(2)} + b^{(3)}
 \end{aligned}$$

$$a^{(2)} = \sigma(z^{(2)})$$

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

$$\frac{\partial L}{\partial w^{(2)}} = \underbrace{-(y^{(i)} - a^{(3)})}_{(1 \times 1)} \underbrace{(w^{(3)})^T}_{(2 \times 1)} \underbrace{a^{(2)} (1-a^{(2)})}_{(2 \times 1)} \underbrace{(a^{(1)})^T}_{(1 \times 3)}$$

$$\frac{\partial L}{\partial a^{(3)}} \times \frac{\partial a^{(3)}}{\partial z^{(3)}} \quad \frac{\partial z^{(3)}}{\partial a^{(2)}} \quad \frac{\partial a^{(2)}}{\partial z^{(2)}} \quad \frac{\partial z^{(2)}}{\partial w^{(2)}}$$

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

$$w^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \\ w_{21}^{(2)} & w_{22}^{(2)} & w_{23}^{(2)} \end{bmatrix}$$

Subject:

$$\frac{\partial z^{(2)}}{\partial w^{(2)}} = \begin{bmatrix} \frac{\partial z^{(2)}}{\partial w_{11}^{(2)}} & \frac{\partial z^{(2)}}{\partial w_{12}^{(2)}} & \frac{\partial z^{(2)}}{\partial w_{13}^{(2)}} \\ \frac{\partial z^{(2)}}{\partial w_{21}^{(2)}} & \frac{\partial z^{(2)}}{\partial w_{22}^{(2)}} & \frac{\partial z^{(2)}}{\partial w_{23}^{(2)}} \end{bmatrix}$$

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

$$= \begin{bmatrix} w_{11}^{(2)} a_1^{(1)} + w_{12}^{(2)} a_2^{(1)} + w_{13}^{(2)} a_3^{(1)} + b^{(2)} \\ w_{21}^{(2)} a_1^{(1)} + w_{22}^{(2)} a_2^{(1)} + w_{23}^{(2)} a_3^{(1)} \end{bmatrix}$$

$$\frac{\partial L}{\partial w^{(2)}} = -(y^{(i)} - a^{(3)}) (w^{(3)})^T a^{(2)} (1 - a^{(2)}) (a^{(1)})^T$$

$$= (w^{(3)})^T \cdot a^{(2)} (1 - a^{(2)}) \left(-(y^{(i)} - a^{(3)}) (a^{(1)})^T \right)$$

Component wise product $a^{(2)} = \sigma(z^{(2)})$

$$\begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} = \sigma \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \\ \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \end{bmatrix} = \begin{bmatrix} a_1^{(2)} (1 - a_1^{(2)}) \\ a_2^{(2)} (1 - a_2^{(2)}) \end{bmatrix}$$

$$\text{Component-wise product : } \begin{bmatrix} a \\ b \end{bmatrix}_{2 \times 1} \begin{bmatrix} c \\ d \end{bmatrix}_{2 \times 1} = \begin{bmatrix} a.c \\ b.d \end{bmatrix}_{2 \times 1}$$

$$\frac{\partial L}{\partial w^{(2)}} = \underbrace{(w^{(3)})^T \cdot a^{(2)} (1 - a^{(2)})}_{(2 \times 1)} \underbrace{\left(\frac{-(y^{(i)} - a^{(3)}) (a^{(1)})^T}{(1 \times 1) (1 \times 3)} \right)}_{(1 \times 3)}$$

★ Support vector Machines (SVMs) :-

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

$$\theta^T x >> 0 \longrightarrow y = 1$$

$$\theta^T x << 0 \longrightarrow y = 0$$

Functional Margin —

In SVM, such insights fall under the notion of functional margins.

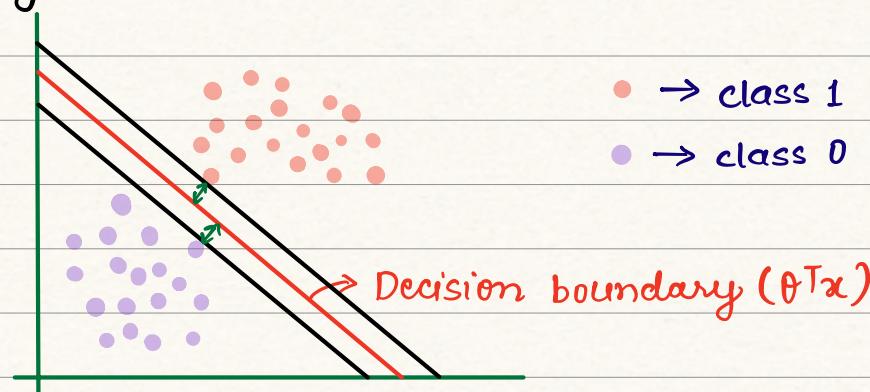
To avoid $g(\theta^T x) \approx 0.5$

We will train the parameters (θ) using gradient descent algorithm to get ;

(i) $g(\theta^T x) \gg 0.5$ for $y = 1$

(ii) $g(\theta^T x) << 0.5$ for $y = 0$

Geometric Margins —



Notations :—

$y \in \{-1, 1\}$
(SVM)

$y \in \{0, 1\}$
(Logistic Regression)

$h_{w,b}(x) = g(w^T x + b)$
(SVM)

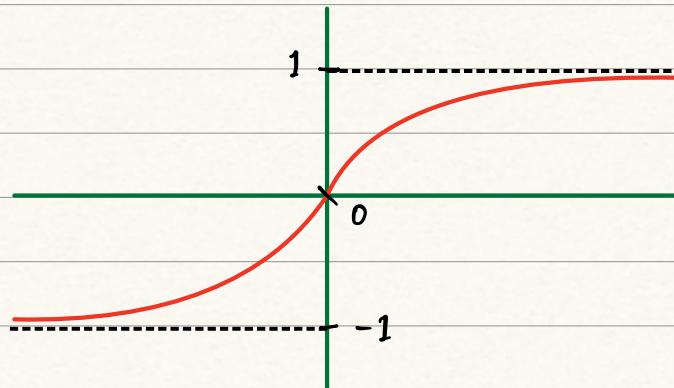
$h_\theta(x) = g(\theta^T x)$
(Logistic Regression)

Subject:

/ /

$$g(w^T x + b) = 1 \quad \text{if } w^T x + b \geq 0$$

$$g(w^T x + b) = -1 \quad \text{if } w^T x + b < 0$$



* Functional Margin in context of SVM :—

Suppose I have a training example $(x^{(i)}, y^{(i)})$

↙ ↘ $\{ -1, 1 \}$

Input feature vector
for i th training example

Functional margin of (w, b) w.r.t this training example is ;

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b) \quad [y^{(i)} = \{-1, +1\}]$$

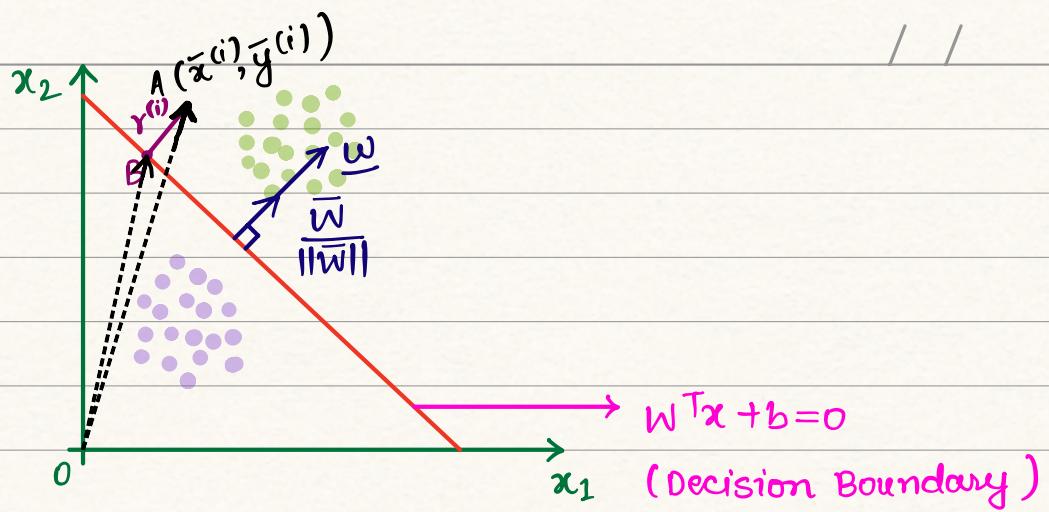
If $\hat{\gamma}^{(i)} > 0 \rightarrow$ we are not making any training error.

$$h_{w,b}(x^{(i)}) = y^{(i)}$$

$$S = [x^{(i)}, y^{(i)}] ; \quad i = 1, 2, \dots, n$$

$$\hat{\gamma} = \min_{1, 2, \dots, n} \hat{\gamma}^{(i)}$$

* Geometric Margin in context of SVM :—



$$\vec{A} - \vec{B} = \gamma^{(i)} \frac{\vec{w}}{\|\vec{w}\|}$$

$$\vec{B} = \vec{A} - \gamma^{(i)} \frac{\vec{w}}{\|\vec{w}\|}$$

For point A ;

$$\vec{B} = \vec{x}^{(i)} - \gamma^{(i)} \frac{\vec{w}}{\|\vec{w}\|}$$

For point B ;

$$w^T \vec{B} + b = 0$$

$$w^T \left(\vec{x}^{(i)} - \gamma^{(i)} \frac{\vec{w}}{\|\vec{w}\|} \right) + b = 0$$

$$\Rightarrow w^T \vec{x}^{(i)} - \gamma^{(i)} \frac{w^T w}{\|\vec{w}\|} + b = 0$$

$$\{ w^T w = \|w\|^2 \}$$

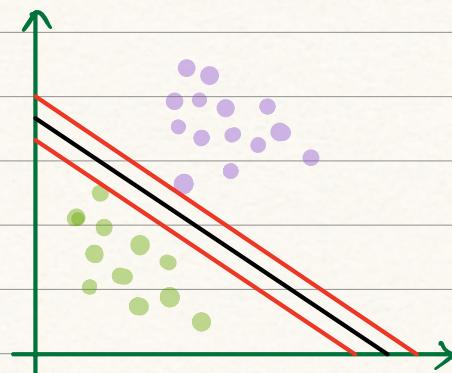
$$\Rightarrow w^T \vec{x}^{(i)} - \gamma^{(i)} \|w\| + b = 0$$

$$\gamma^{(i)} = \frac{w^T \vec{x}^{(i)} + b}{\|w\|}$$

In this case, the Geometric Margin for the $(x^{(i)}, y^{(i)})$ training example is ;

$$\gamma^{(i)} = y^{(i)} \left(\frac{w^T \vec{x}^{(i)}}{\|w\|} + \frac{b}{\|w\|} \right)$$

if $\|w\| = 1$;
 $\rightarrow \text{Geometric Margin} \equiv \text{Functional Margin}$



Hard SVM

Clear separation
between two datasets

optimization problem :

$$\max_{w, b} \gamma \quad (\text{Geometric Margin})$$

$$\text{subject to } y^{(i)} (w^T x^{(i)} + b) \geq \gamma \quad \{i=1, 2, \dots, n\}$$

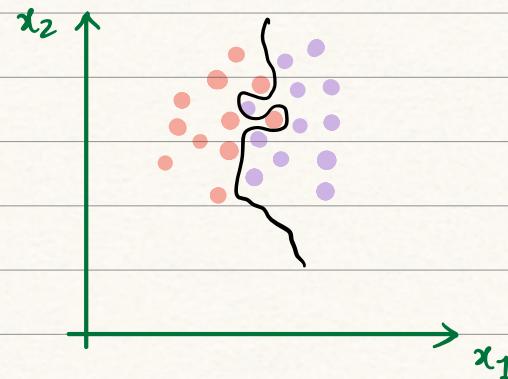
$$\|w\| = 1$$

Q.) why should I go beyond the hard SVMs?

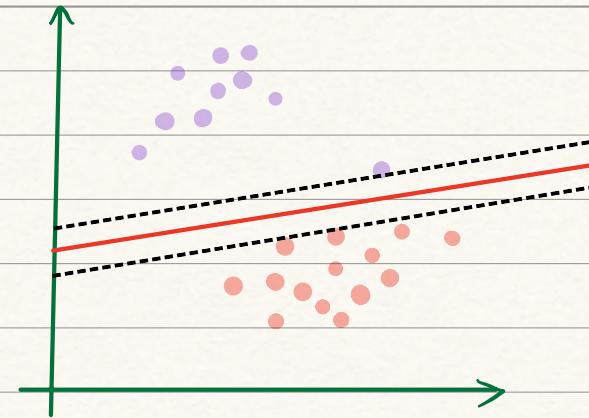
→ We should opt for soft SVMs.

→ We should relax some constraints.

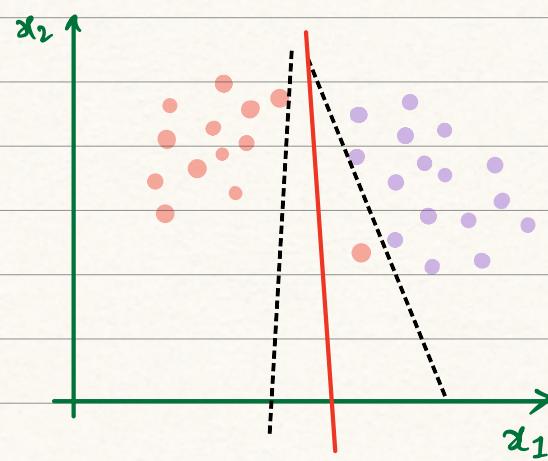
(i) In general, the dataset in realistic applications is non-separable.



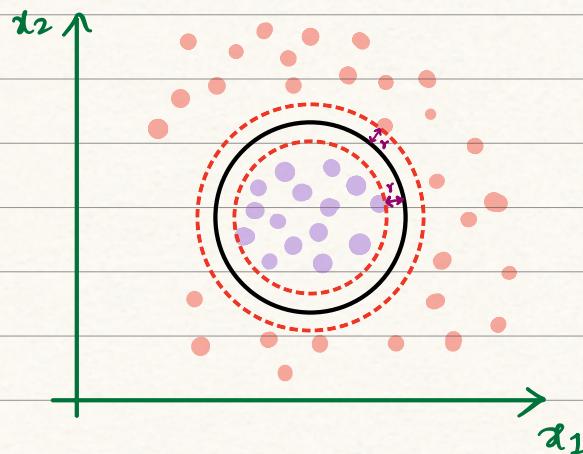
(ii) Hard SVMs would not work when there are some outliers in the dataset.



Geometric margin
is minimized due
to outlier.



(soft Margin)



Transformation
using kernels
(Combine SVMs
with kernels)

■ Unsupervised Learning Algorithms :

★ Clustering :—

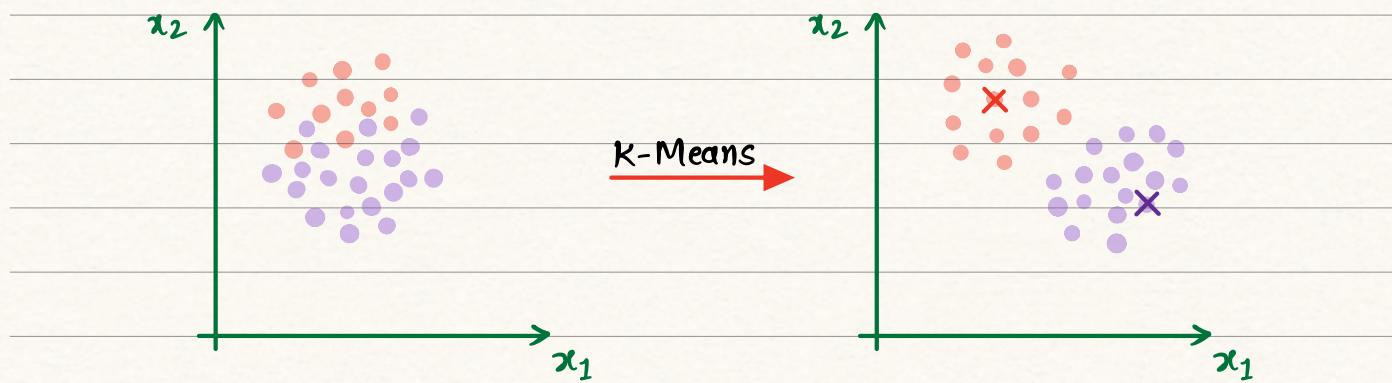
Suppose we have a training set ;

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

and we want to group the data into a few classes.

Here, $x^{(i)} \in \mathbb{R}^d$ and there is no labelling.

So this is unsupervised learning problem.



* K-Means clustering Algorithm :-

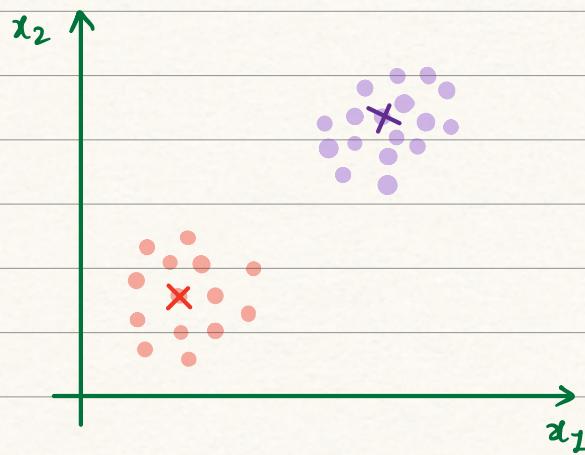
1.) Initialize the cluster centroids $\bar{\mu}_1, \bar{\mu}_2, \dots, \bar{\mu}_k \in \mathbb{R}^d$ randomly.

2.) Repeat until convergence for every i , set ;

$$c^{(i)} := \arg \min_j \| \bar{x}^{(i)} - \bar{\mu}_j \|^2$$

For every j (clusters), set ;

$$\mu_j := \frac{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = j\} \cdot x^{(i)}}{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = j\}}$$



★ Expectation-Maximization (EM) Algorithm for Mixture of Gaussians —

$\{\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(n)}\} \rightarrow$ Data do not have any labels

$$\underbrace{(\{\underline{x}^{(1)}, \underline{x}^{(2)}, \underline{x}^{(3)}\})}_{(\underline{\mu}_1, \Sigma_1)}$$

$$\underbrace{(\{\underline{x}^{(4)}, \underline{x}^{(5)}, \underline{x}^{(6)}\})}_{(\underline{\mu}_2, \Sigma_2)}$$

We first model the input data by specifying a Joint distribution ;

$$f(\underline{x}^{(i)}, \underline{z}^{(i)}) = f(\underline{x}^{(i)} | \underline{z}^{(i)}) P(\underline{z}^{(i)})$$

$\underline{z}^{(i)} \sim \text{Multinomial distribution } (\Phi)$

$$\phi_j \geq 0, \sum_{j=1}^k \phi_j = 1, \phi_j = P(z^{(i)}=j)$$

$$(\underline{x}^{(i)} | z^{(i)}=j) \sim \mathcal{N}(\underline{\mu}_j, \Sigma_j)$$

Q.) How $\underline{x}^{(i)}$ was generated ?

$$i = \{1, 2, 3, \dots, k\} \quad (\text{k Gaussians})$$

Randomly choose Gaussian from 'i'

$\underline{x}^{(i)}$ dataset coming Gaussian 'i'

$$\underbrace{\{\underline{x}^{(1)}, \underline{x}^{(2)}, \underline{x}^{(3)}, \underline{x}^{(4)}, \underline{x}^{(5)}, \underline{x}^{(6)}\}}$$

$$(\underline{\mu}_1, \Sigma_1) \quad (\underline{\mu}_2, \Sigma_2) \quad (\underline{\mu}_3, \Sigma_3)$$

Subject:

$$f(\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(n)}) = f(\underline{x}^{(1)}) \cdot f(\underline{x}^{(2)}) \cdots f(\underline{x}^{(n)})$$

$$= \prod_{i=1}^n f(\underline{x}^{(i)})$$

$$L(\underline{\mu}, \Sigma, \phi) = \prod_{i=1}^n f(\underline{x}^{(i)}, \underline{\mu}, \Sigma, \phi)$$

↑ likelihood

$$\log(L(\underline{\mu}, \Sigma, \phi)) = \sum_{i=1}^n \log f(\underline{x}^{(i)}, \underline{\mu}, \Sigma, \phi)$$

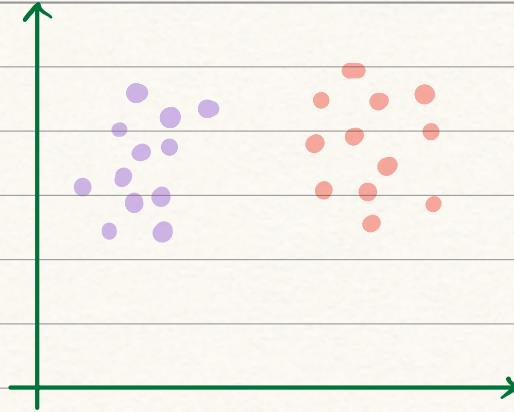
$$= \sum_{i=1}^n \log \sum_{\underline{z}^{(i)}=1}^K f(\underline{x}^{(i)} | \underline{z}^{(i)}, \underline{\mu}, \Sigma) p(\underline{z}^{(i)}, \phi)$$

If we know $\underline{x}^{(i)}$ is coming from $\underline{\mu}^{(i)}$ Gaussian.
(labelled data)

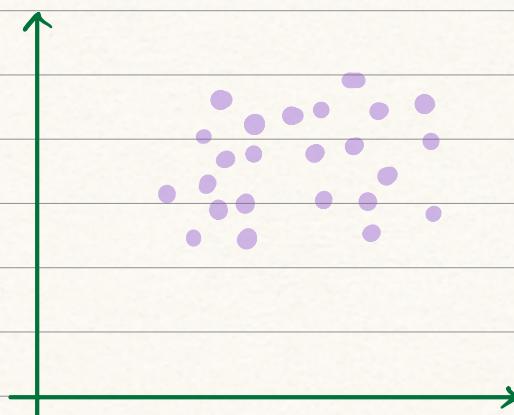
$$\phi_j = \frac{\sum_{i=1}^n \mathbb{1}\{\underline{x}^{(i)} = j\}}{n}$$

$$\underline{\mu}_j = \frac{\sum_{i=1}^n \mathbb{1}\{\underline{x}^{(i)} = j\} \cdot \underline{x}^{(i)}}{\sum_{i=1}^n \mathbb{1}\{\underline{x}^{(i)} = j\}}$$

$$\Sigma_j = \frac{\sum_{i=1}^n \mathbb{1}\{\underline{x}^{(i)} = j\} (\underline{x}^{(i)} - \underline{\mu}_j)(\underline{x}^{(i)} - \underline{\mu}_j)^T}{\sum_{i=1}^n \mathbb{1}\{\underline{x}^{(i)} = j\}}$$



$K = 2$ Gaussians
(Labelled data)



$K = 2$ Gaussians
(Unlabelled data)

So for unlabelled data we will use EM Algorithm for mixture of Gaussians.

It has two steps : —

1.) E-step : It tries to guess the values of $z^{(i)}$'s.

2.) M-step : It updates the parameters.

Repeat until it converges {
(E-step). For each i, j , set
 $w_j^{(i)} := p(z^{(i)}=j | x^{(i)}; \phi, \Sigma, \mu)$

Weights
(0, 1)
Soft clustering

Clusters coming from two
or more Gaussians.

what is the probability that the
ith training example is drawn
from jth Gaussian once I
observed it.

Subject:

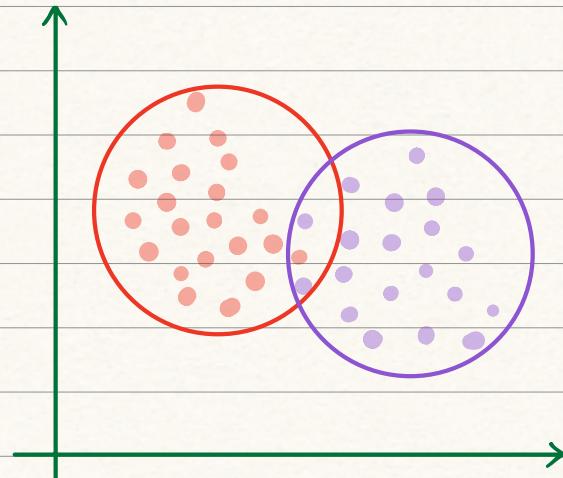
/ /

(M-step);

$$\phi_j := \frac{\sum_{i=1}^n w_j^{(i)}}{n}$$

$$\underline{\mu}_j := \frac{\sum_{i=1}^n w_j^{(i)} \underline{x}^{(i)}}{\sum_{i=1}^n w_j^{(i)}}$$

$$\Sigma_j := \frac{\sum_{i=1}^n w_j^{(i)} (\underline{x}^{(i)} - \underline{\mu}_j)(\underline{x}^{(i)} - \underline{\mu}_j)^T}{\sum_{i=1}^n w_j^{(i)}}$$



$$p(\underline{x}^{(i)} = j | \underline{x}^{(i)}, \phi, \mu, \Sigma) = \frac{f(\underline{x}^{(i)} | z^{(i)} = j; \mu, \Sigma) p(z^{(i)} = j, \phi)}{\sum_{k=1}^K f(\underline{x}^{(i)} | z^{(i)} = k; \mu, \Sigma) p(z^{(i)} = k, \phi)}$$