

# Machine Learning

Definition:

- ① ML is a field of study that gives computers the ability to learn without being explicitly programmed.

[Arthur Samuel, 1959].

- ② A computer program is said to learn from experience E w.r.t. some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.

## ML Algorithms

- Supervised learning
- Unsupervised learning

Others: Reinforcement learning, recommender systems.

- ① Supervised learning :- "right answers" given. labelled output.

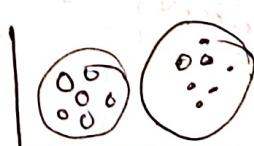
Regression: Predict continuous valued output.

Classification: Discrete valued output (0 or 1)

- ② Unsupervised learning :- no label or same label of all data.

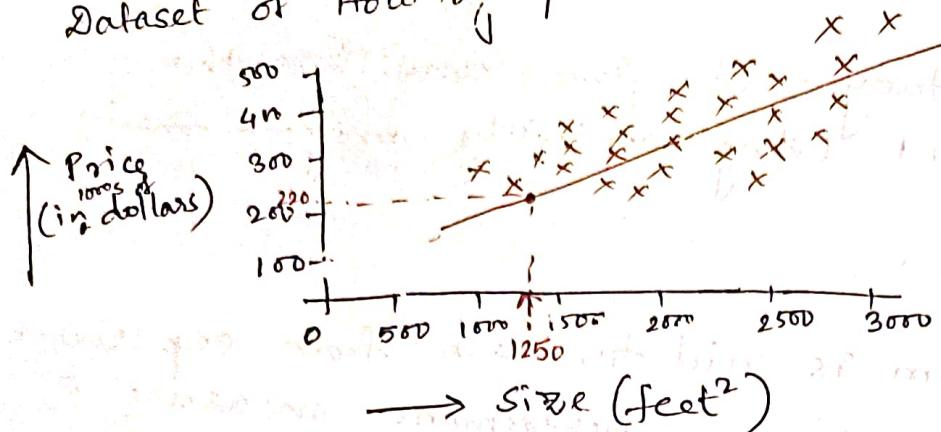
### clustering:

- used to organize large computer clusters
- Application on social network analysis.
- Market segmentation.



## ⑦ Linear Regression with one variable $\Rightarrow$

Dataset of Housing prices.



for a house of size 1250 feet<sup>2</sup>, how many prices will be valued? considered?  
- may be 220 dollars.

A number of houses that were different sizes and sold for a range of different prices

### Supervised

Given the "right answer" for each example in the data

### Regression

Predict real-valued output.

Most commonly used / supervised algorithm is

classification, that gives discrete valued output

Training set of housing prices.

	Size in feet <sup>2</sup> ( $x$ )	Price (\$) in 1000's ( $y$ )
	2104	460
	1416	232
	1534	315
	852	178
	...	...

Let  $m$  be the number of training examples.  
 $x$ : i/p variables/features ,  $y$ : "o/p variable / target variable".

Here, in our example,  $m = 47$ .

$(x, y)$  - one training example.

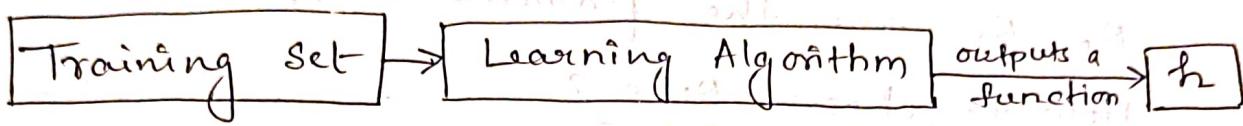
$(x^{(i)}, y^{(i)})$  -  $i^{th}$  training example.

i.e.,  $x^{(1)} = 2104$

$x^{(2)} = 1416$

$y^{(1)} = 460$ .

$h$  stands for hypothesis.



size of house  $\Rightarrow [h] \Rightarrow$  estimated price.

So,  $h$  is a function, (called hypothesis), that maps from  $x$ 's to  $y$ 's.

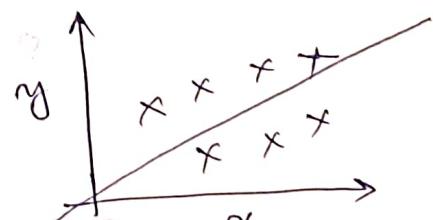
» How do we represent  $h$ ?

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

$h_{\theta}(x)$  is also written as,  $h(x)$ .

It means that, we are going to predict that  $y$  is a linear function of  $x$ .

- Sometimes, linear function does not work. Non linear has to be incorporated. But since linear case is the first (simple) building block, we start with this example of fitting a linear functions.



This model is called linear regression.  
 $h(x) = \theta_0 + \theta_1 x$  is linear regression with one variable. Predicting all the prices as functions of one variable  $x$ .

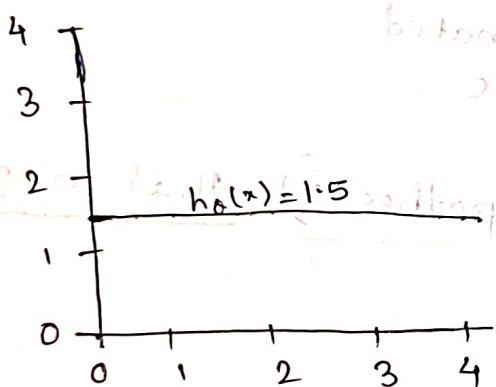
Also called, univariate linear regression.

## ④ Cost Function

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

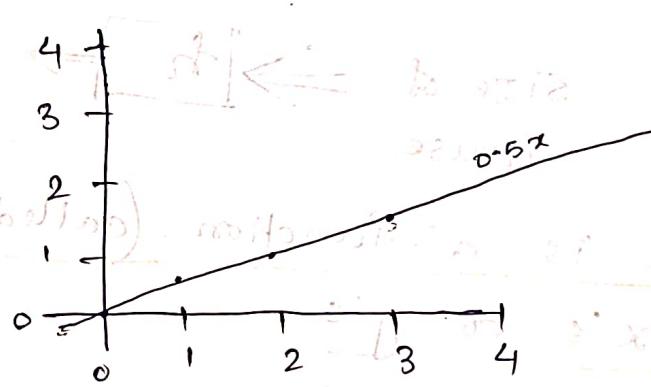
$\theta_0, \theta_1$  are the parameters of the model.

How to choose  $\theta_i$ 's?

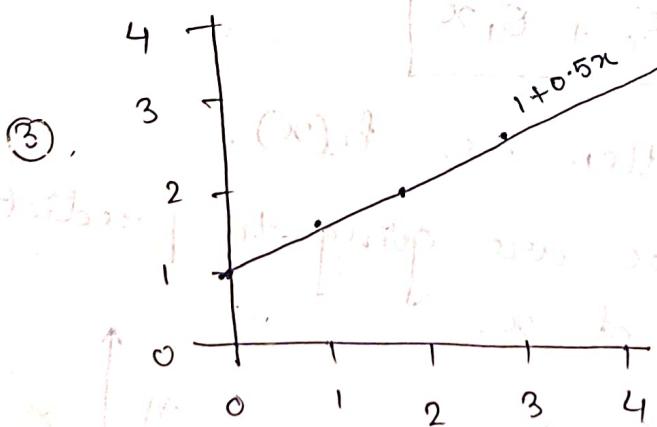


$$\theta_0 = 1.5$$

$$\theta_1 = 0$$

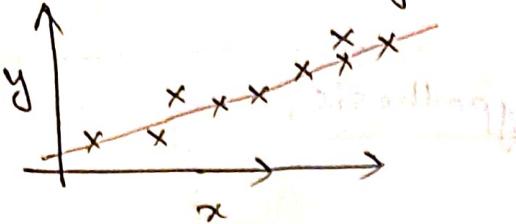


$$\begin{cases} \theta_0 = 0 \\ \theta_1 = 0.5 \end{cases}$$



$$\begin{cases} \theta_0 = 1 \\ \theta_1 = 0.5 \end{cases}$$

$$\textcircled{1} \quad h_{\theta}(x) = 1.5 \quad , \quad \textcircled{2} \quad h_{\theta}(x) = 0.5x \quad , \quad \textcircled{3} \quad h_{\theta}(x) = 1 + 0.5x$$

» In linear regression, we have a training set.   

 we want to come up with values for the parameters  $\theta_0$  &  $\theta_1$ , so that the straight line we get, corresponds to a straight line that somehow fits the data well (like in the above figure).

So, how to get such  $\theta_0$  &  $\theta_1$ , that have good fit to data?

**Idea:** Choose  $\theta_0$  &  $\theta_1$  so that  $h_\theta(x)$  is close to  $y$  for our training examples  $(x, y)$ . ( $\theta = \theta_0 + \theta_1 x$ )  
 ( $y$  is the actual price),  $h$  is the hypothesis.  
 $y$  &  $h$  should be close to each other.

So, in linear regression, we want to solve a minimization problem.

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

Where,  $h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$ ,

This is the overall objective function for linear regression (deciding  $\theta_0, \theta_1$ ).

So, Cost function is,

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

$$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$$

cost function (also called, squared error function).

The squared error func<sup>n</sup> is most commonly used cost function.

However, there are other alternative cost function as well.

Let us now use a simplified hypothesis,

$$h_0(x) = \theta_1 x, \quad \theta_0 = 0$$

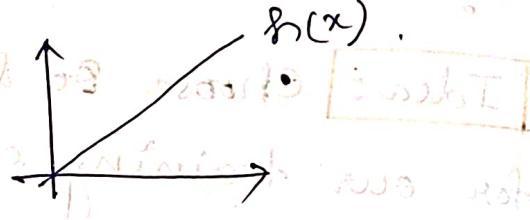
$$\therefore J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

Goal:

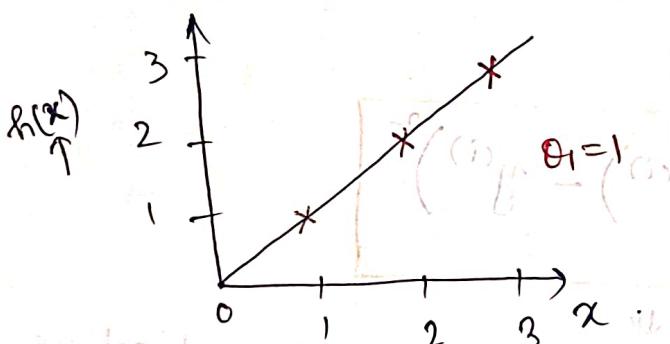
minimize  $J(\theta_1)$

at  $\theta_1$  of  $h_0(x)$  is zero as  $\theta_1$  is zero

When,  $\theta_0 = 0$ , this means  $\rightarrow$



$h_0(x)$  is a function of the parameter  $\theta_1$  (for fixed  $\theta_0$ , this is a function of  $x$ )



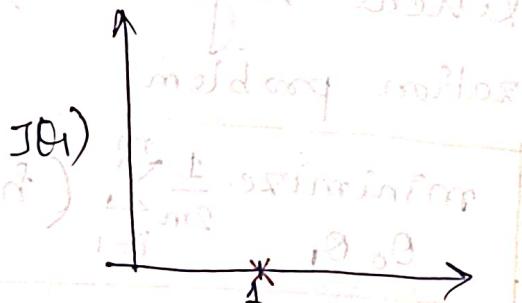
$$\therefore J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

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

Here,  $x = y$  on the above graph.

$$\text{So, } J(\theta_1) = \frac{1}{2m} \times (\theta_1^2 + \theta_1^2 + \dots + \theta_1^2)$$

$$\text{i.e., } J(\theta_1) = 0, \quad \text{Here } m = 3.$$

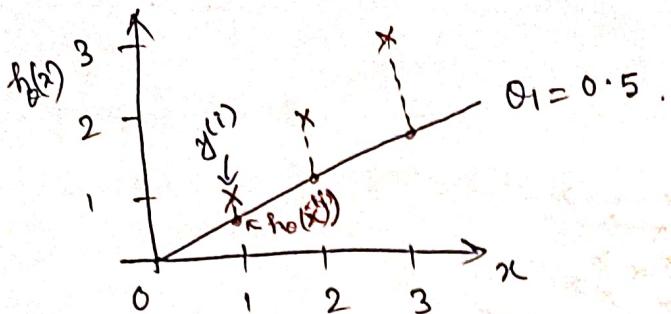


$$\theta_1 = 0$$

$$J(0) = 0$$

$$\frac{\partial J}{\partial \theta_1} = (0, 0)$$

$$\theta_1 = 0.5$$



$$\therefore J(0.5)$$

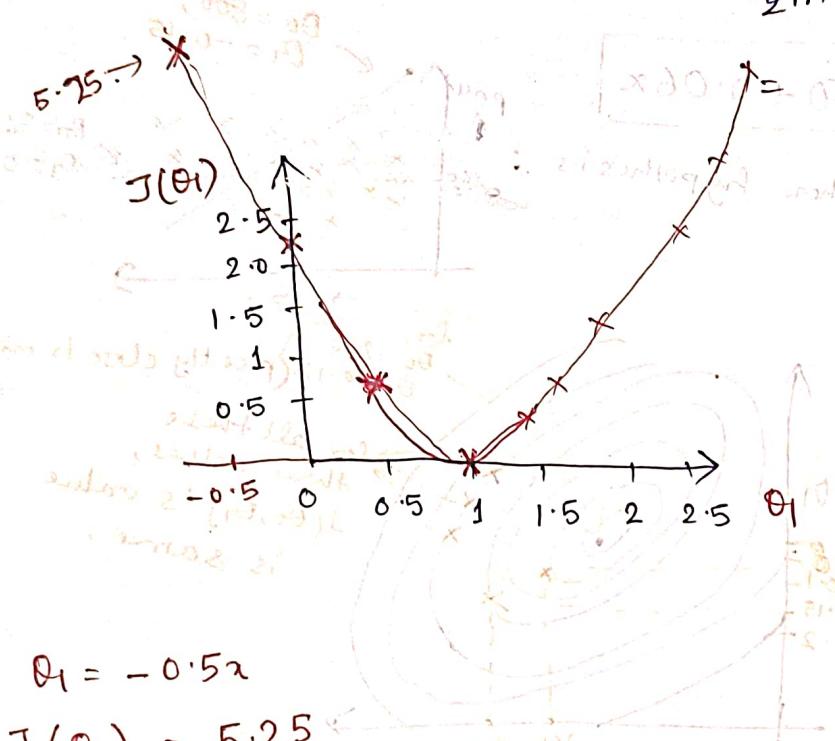
$$= \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} \times \sum_{i=1}^m (0.5x^{(i)} - y^{(i)})^2$$

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

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

$$\& h(x) = 0.5x$$



$$\theta_1 = -0.52$$

$$J(\theta_1) = 5.25$$

$$\theta_1 = 1, J(\theta_1) = 0$$

$$\theta_1 = 0.5, J(\theta_1) = 0.58$$

$$\theta_1 = 0, J(\theta_1) =$$

$$J(\theta_1) = \frac{1}{6} [1^2 + 2^2 + 3^2]$$

$$= \frac{1}{6} \times 14 = \frac{14}{6} = 2.3$$

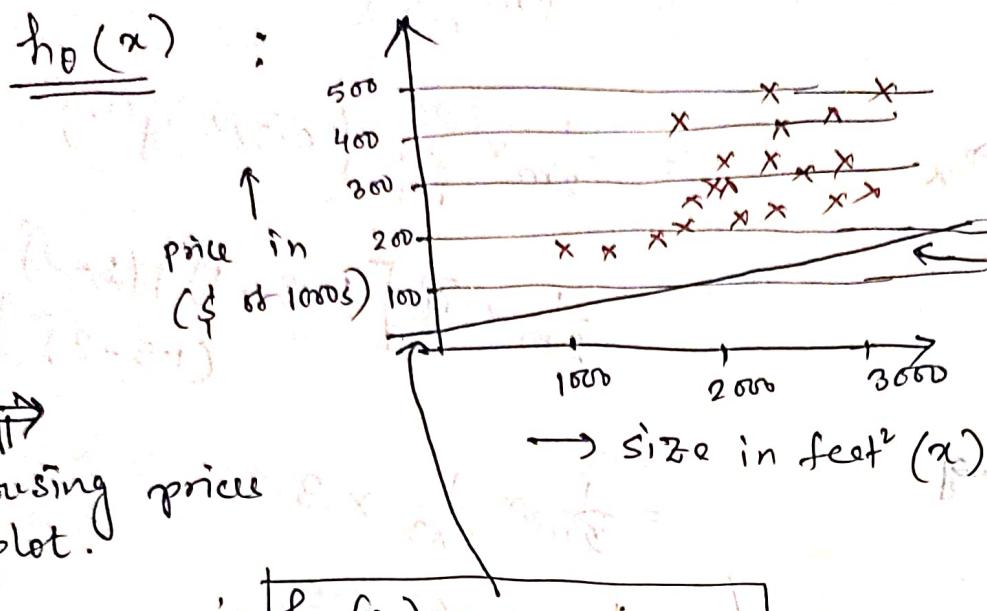
Our objective is the optimization objective for our learning algorithm, so we want to choose the value of  $\theta_1$  so that it minimises  $J(\theta_1)$ . This is the

minimize  $J(\theta_1)$   $\leftarrow$  This is our objective function of linear regression.

Here, in the above example, the value that minimizes  $J(\theta_1)$  is,  $\theta_1 = 1$ . That is the possible best line fit through our data.

>> Let us now keep both the parameters  $\theta_0$  &  $\theta_1$

$$\underline{h_0(x)}$$

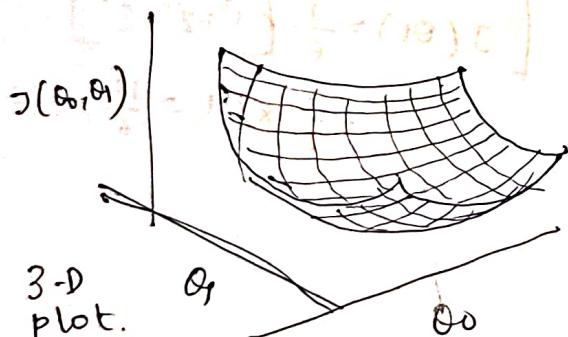


housing price plot.

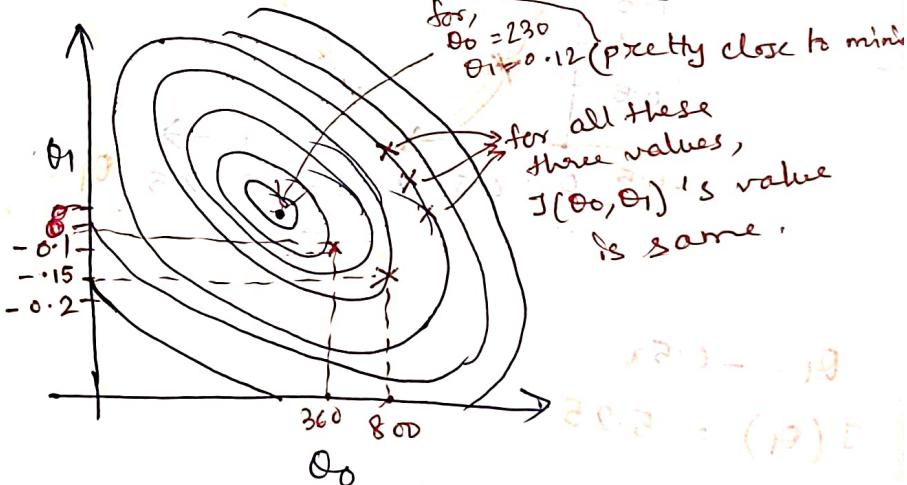
$$\therefore \underline{h_0(x) = 50 + 0.06x}$$

Now,  $J(\theta_0, \theta_1)$ ,  $\lambda = 10$

$$\underline{J(\theta_0, \theta_1)} : \theta_0 = 10$$



Contour plot



\* Each of the ellipses take ~~sorry~~ shows a set of points that takes on the same value of  $\theta_0$  &  $\theta_1$  for  $J(\theta_0, \theta_1)$ .

- ① Contour figure is a way (or a more convenient way) to ~~change~~ visualize my function  $J$ .

① Example : If  $\theta_0 = 800$  &  $\theta_1 = -0.15$  is not a very good fit to our example (as shown in the previous page). The value of  $J(\theta_0, \theta_1)$  is pretty far from the minimum (the centre) (as shown in the previous page). i.e., pretty high cost of  $J(\theta_0, \theta_1)$  we obtain.

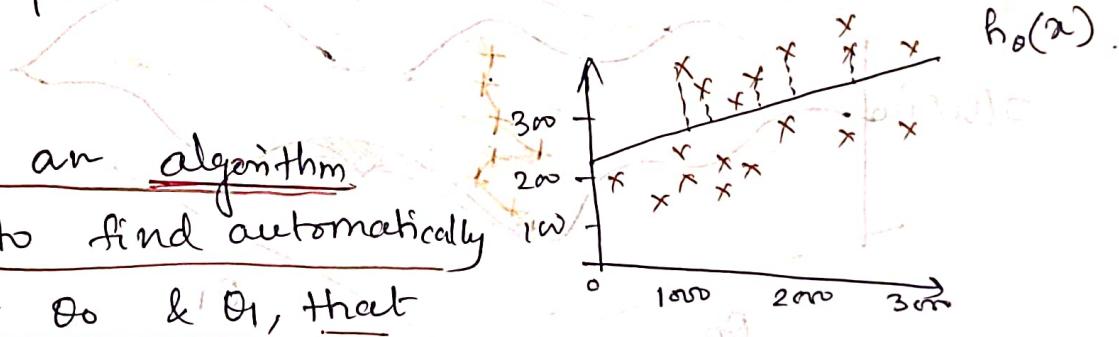
Again, if  $\theta_0 = 360$ ,  $\theta_1 = 0$ . (closer)

Again,  $\theta_0 = 500$ ,  $\theta_1 =$

then  $\theta_0 = 230$ ,  $\theta_1 = 0.12$ ,  $\rightarrow J(\theta_0, \theta_1)$  is pretty close to minimum.

[smallest & centered ellipse have same height, next ellipse has height greater than the central one, ... so on.]

[the smallest ellipse also has minimum value].



So, we need an algorithm or software to find automatically the values of  $\theta_0$  &  $\theta_1$ , that minimizes cost function  $J$ .

② The Algorithm : Gradient descent algorithm.

It is the most general algorithm to find out values of  $\theta_0$ ,  $\theta_1$  that minimizes  $J(\theta_0, \theta_1)$ . It is not only used in linear regression, but also is used all over the places of machine learning and helps in minimizing others function as well apart from the cost function.

- ① We have some function  $J(\theta_0, \theta_1)$ .  
 We want  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ .
- Outline:
- ① start with some  $(\theta_0, \theta_1)$
  - ② keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$   
 until we hopefully end up at a minimum!

[Gradient descent is actually applied to more general function, like  $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$  and

want  $\boxed{\min_{\theta_0, \theta_1, \dots, \theta_n} J(\theta_0, \theta_1, \dots, \theta_n)}$

Say, initialize  $\theta_0 = 0, \theta_1 = 0$  : starting point.



- imagine you are standing on a hill.
- you want to go downhill as quickly as possible.
- you take little baby steps in some direction.
- what direction do I take that little baby step in?
- you choose a direction and go down and down until you get the local minimum.

- similarly, by choosing another direction, you can reach to another local minimum (which is very different from the previous one).

### Algorithm :

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0 \text{ and } j=1)$$

}

- $\alpha$  is a number that is called learning rate.
- $\alpha$  denotes how big step we take downhill with creating descent.
- If  $\alpha$  is very large, then it corresponds to a very aggressive gradient descent procedure where we are trying to take huge step downhill.
- Simultaneously update  $\theta_0, \theta_1$ .

i.e.,

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

### Incorrect:

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}_0$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 := \text{temp}_1$$

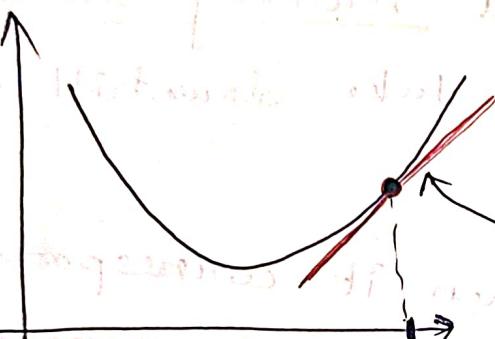
④ Example:  $\theta_0 = 1, \theta_1 = 2$   
 $\theta_j := \theta_j + \sqrt{\theta_0 \theta_1}$  (for  $j=0 & 1$ )  
 simultaneously update  $\theta_0$  &  $\theta_1$ . what are the resulting values of  $\theta_0$  &  $\theta_1$ ?

Ans:

$$\theta_0 = \theta_1 + \sqrt{2} \quad (3, 3) \rightarrow \theta_0 = 3 + \sqrt{2} \approx 4.236$$

$$\theta_1 = 2 + \sqrt{2}$$

⑤



For simplicity let's assume, we have only  $\theta_1$ . Let's start with initial value  $\theta_1$ .

derivative term denotes derivative of  $J(\theta_1)$  w.r.t.  $\theta_1$ . Here, slope is positive.

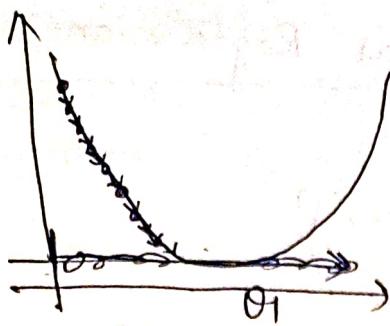
$$\text{So, } \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1) \rightarrow \text{this term is positive } (z, o)$$

$$\text{i.e., } \theta_1 := \theta_1 - \alpha. \quad (\text{some positive value})$$

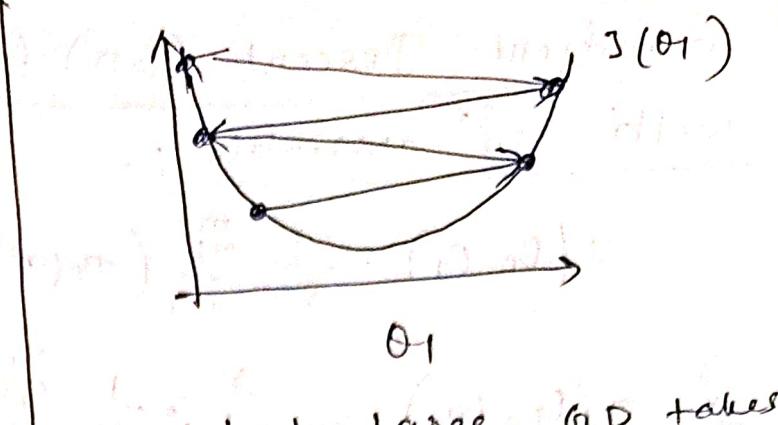
So, new updated  $\theta_1$  will have less value compared to the previous  $\theta_1$ .

Similarly, if in some point, slope is negative then, updated  $\theta_1$  will have more increased value than the previous one.

- ⑥ If  $\alpha$  is too small, gradient descent can be slow.
- ⑦ If  $\alpha$  is too large, gradient descent can overshoot minimum, hence may fail to converge or diverge.



when  $\alpha$  is too small

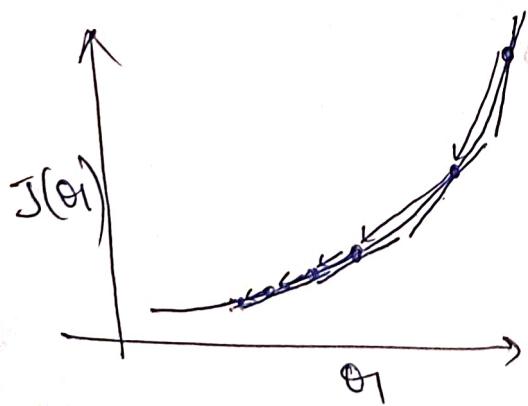


$\alpha$  is too large, GD takes huge step. It goes further away from minimum.

Qs:- what if the local minimum is at parameter  $\theta_1$  is already in local minimum?

Ans: slope is horizontal. So, derivative value is 0. and  $\theta_1' = \theta_1$  i.e. value of  $\theta_1$  remains unchanged

② even if the learning rate  $\alpha$  is fixed, GD can converge to a local minimum. As we approach a local minimum, GD will automatically take smaller steps. So, no need to decrease  $\alpha$  over time.



The less the steepness of the slope, the less the value of the derivative term. So, the smaller derivative term helps in taking the smaller steps (or baby steps) even if  $\alpha$  is fixed.

## ① Gradient Descent (GD) for Linear Regression with one variable :-

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

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)})^2$$

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

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

$$= \frac{2}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}). \text{ Differentiation of } \theta_0.$$

$$= \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) = \frac{1}{m} \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}).$$

And, similarly differentiate with  $\theta_1$  and

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

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

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

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

So, now, ~~GD~~

gradient descent algorithm for linear regression is,  
repeat until convergence {

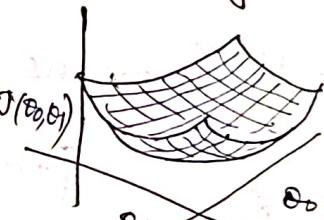
$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\hat{h}_0(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\hat{h}_0(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

} update  $\theta_0$  &  $\theta_1$  simultaneously.

- ① Cost function of linear regression is always going to be a bow shaped function.

The technical term of this function is, convex function.



Therefore, it always converges to the global optimum (there is no other local optimum).

- ② The GD algorithm, ~~in this case, is known machine learning~~, is called Batch Gradient Descent.

"Batch": Each step of gradient descent uses all the training examples.

$$\sum_{i=1}^m$$

everytime / when we <sup>update</sup> calculate

$\theta_0$  or  $\theta_1$ .

### ⑤ House sizes:

$$2104$$

$$1416$$

$$1534$$

$$852$$

my matrix is,

$$\begin{bmatrix} 2104 \\ 1416 \\ 1534 \\ 852 \end{bmatrix}$$

and vector is

$$\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$$

$$\therefore \begin{bmatrix} 1 & -40 \\ 1416 & 0.25 \end{bmatrix} =$$

$$\begin{bmatrix} -40 + 0.25 \times 2104 \\ -40 + 0.25 \times 1416 \\ -40 + 0.25 \times 1534 \\ -40 + 0.25 \times 852 \end{bmatrix}$$

$$f_0(2104)$$

$$\therefore \text{Prediction} = \text{DataMatrix} \times \text{Parameters}$$

gives all the predictive prices.

$$\Rightarrow \begin{bmatrix} 1 & -\frac{3}{2} \\ 1416 & 1 \\ 1534 & 1 \\ 852 & 1 \end{bmatrix}_{2 \times 3} \times \begin{bmatrix} 1 & 3 \\ 0 & 1 \\ 5 & 2 \end{bmatrix}_{3 \times 2} =$$

$$\begin{bmatrix} 1 \times 1 + 3 \times 0 + 2 \times 5 \\ 1 \times 3 + 3 \times 1 + 2 \times 2 \\ 4 \times 3 + 0 \times 1 + 1 \times 2 \end{bmatrix}_{3 \times 2}$$

$$= \begin{bmatrix} 11 & 10 \\ 19 & 14 \end{bmatrix}_{3 \times 2}$$

Ans.

### ⑥ House sizes:

$$2104$$

$$1416$$

$$1534$$

$$852$$

$$1. f_0(x) = -40 + 0.25x$$

$$2. f_0(x) = 200 + 0.1x$$

$$3. f_0(x) = -150 + 0.4x$$

Have 3 competing hypotheses:

$$\therefore \begin{bmatrix} 1 & -40 \\ 1416 & 0.25 \end{bmatrix} \times \begin{bmatrix} 200 & -150 \\ 0.1 & 0.4 \end{bmatrix}_{2 \times 2} =$$

$$\begin{bmatrix} 486 & 410 \\ 314 & 342 \\ 344 & 353 \\ 285 & 191 \end{bmatrix}_{4 \times 2} = \text{predicted housing prices.}$$

$$m = 4, (4 \text{ houses})$$

∴ 12 predictions we get.

⑦ Matrix Inverse: If  $A$  is an  $m \times m$  square matrix and if it has an inverse, then if  $A^{-1}$  is an  $m \times m$  matrix, and if it has an inverse, then  $A(A^{-1}) = A^{-1}A = I$ .

⇒ Intuitively, if a matrix does not have any inverse, then it can be thought of as, it is too close to 0 in some sense.

⇒ Matrices that don't have an inverse are "singular" or "degenerate" matrix.

## ① Linear Regression with multiple variables :-

For one variable : size (feet) | price (\$1000)

2104	460
1416	292
1534	315
852	178

So, our hypothesis was,  $h_{\theta}(x) = \theta_0 + \theta_1 x$ .

Multiple features (variables) :-			
size (feet) ( $x_1$ )	# bedrooms ( $x_2$ )	# floors ( $x_3$ )	Age of home (years) ( $x_4$ )
2104	5	1	45
1416	3	2	30
1534	3	2	30
852	2	1	15

Price (\$1000)

460
292
315
178

Notation:

$x^{(i)}$  = # features (here,  $m=4$ )

Parameters  $\theta_0, \theta_1, \dots, \theta_n$  ( $n+1$  dimensional vector  $\theta$ ).

Cost function :

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

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

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}, \quad \text{so, } x_3 = 2.$$

Hypothesis :-

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

$$\therefore \boxed{\begin{aligned} h_{\theta}(x) &= \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \\ \theta &= \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}, \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \end{aligned}} \Rightarrow \text{Multivariate Linear regression.}$$

For convenience, define  $x_0 = 1$ .

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

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

(simultaneously update for every  $j=0, \dots, n$ )

Gradient Descent :

$$\text{Repeat } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$$

$$\text{for all } j = 0, \dots, n$$

## New Algorithm for $m \geq 1$

Repeat

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

(simultaneously update  $\theta_j$  for  $j = 0, 1, \dots, d$ )

» Let us consider 2 features:

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

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

$$\theta_2 := \theta_2 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\phi_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

» It is important to choose features carefully

## Housing price prediction

$$\theta(x) = \theta_0 + \theta_1 x_{\text{frontage}} + \theta_2 x_{\text{depth}}$$

so, two features are there  $\rightarrow x_1$  &  $x_2$ .  
It is not always necessary to use the features  
more features

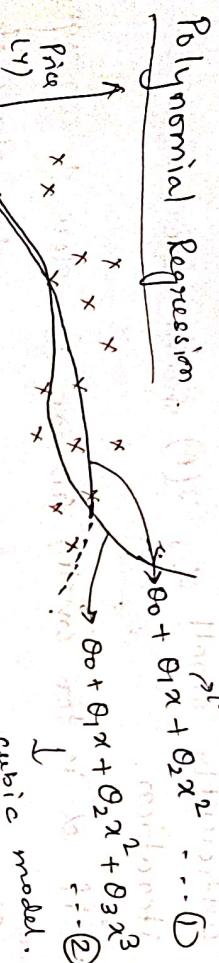
that are given. You can use these  
by yourself. You can also  
use them as a reference.

Let us create a new

Let us assume that the land area is a rectangle of length  $x$  and width  $y$ . Then the area is  $x \cdot y$ . Now, if we consider the frontage as the length  $x$  and the depth as the width  $y$ , then the area is  $x \cdot y$ . This is the same as  $x \cdot 0.4 + 0.4(x^2)$ . So,  $x = \text{frontage} \times \text{depths}$  is Area (  $x$  is the land area).

Sometimes, we suggest

## ① Polynomial Regression



Note that, eq ① gives the curve that will eventually decrease, thus housing price will decrease.

But, eq ② gives the curve (just a cubic value) that increases eventually as size is added to the eq ①), that increases eventually.

So, it is better to use eq ②.

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

$$= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

$$\rightarrow x_1 = \text{size}$$

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

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

Also,  $\theta_0(x)$  may be chosen as  $\theta_0(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 \sqrt{\text{size}}$ .

Then curve will not go down, but it will go up slowly.

So, we can use polynomial functions as well to fit our data.

» To solve the optimal value of parameter  $\theta$ , we sometimes use normal equation (instead of gradient descent).

⇒ In GD, we solve the optimal value of  $\theta$  by taking many steps, multiple iterations of GD to converge to the global minimum.

⇒ In contrast, normal equation is a method to solve for  $\theta$  analytically. So, instead of running iterative algorithm, we can just solve the optimal value for  $\theta$  all in one go (i.e. in one step).

Intuition: If  $f(\theta) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots$

$$f(\theta) = a\theta^2 + b\theta + c \quad (\text{a quadratic func})$$

The way to minimize a quadratic function is by taking derivatives.

$$\frac{d}{d\theta} (J(\theta)) = \dots \stackrel{\text{set}}{=} 0.$$

solve for  $\theta$ :

⇒ It is the simpler case when  $\theta$  is just a real number.

⇒ But we are interested in the case  $\theta$  is  $n \times 1$  dimensional vector, and

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

$$\therefore \frac{\partial}{\partial \theta_j} J(\theta) = \dots \stackrel{\text{set}}{=} 0. \quad (\text{Do for every } j)$$

Examples:  $m = 4$

$\text{size}(x)$	# bedrooms ( $x_1$ )	# floors ( $x_2$ )	Age (years) ( $x_3$ )	Price (floors) ( $y$ )
$x_0$				
1	2104	5	1	45
1	1416	3	2	40
1	1534	3	2	30
1	852	2	1	36
				178

add an extra column.

$m = 4$  training examples

matrix  $X$  is all the features followed by ones.

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

on dimension vector.

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

predict a hill's weight as a function of his age

height with the model.

the value of  $\theta$  that minimizes the cost function.

( $X$  is also called design matrix).

$$\theta_0 + \theta_1 \cdot \text{age} + \theta_2 \cdot \text{height}$$

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

on dimension vector.

is this wrong? because, the ans. is given as,

$$y = \begin{bmatrix} 16 \\ 28 \\ 28 \\ 20 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2104 & 1416 & 1534 & 852 \\ 5 & 3 & 3 & 2 \\ 1 & 2 & 2 & 1 \\ 45 & 40 & 30 & 36 \end{bmatrix}$$

$\therefore X^T X =$

$$= \begin{bmatrix} 4 & 2104+1416+1534+852 & 13 \\ 2104+1416+1534+852 & (2104 \cdot 5) + (1416 \cdot 3) + (1534 \cdot 3) + (852 \cdot 2) & 6 \\ 5 & (2104 \cdot 1) + (1416 \cdot 2) + (1534 \cdot 2) + (852 \cdot 1) & 178 \\ 13 & (2104 \cdot 5) + (1416 \cdot 3) + (1534 \cdot 3) + (852 \cdot 2) & (2104 \cdot 45) + (1416 \cdot 40) + (1534 \cdot 30) + (852 \cdot 36) \end{bmatrix}$$

let,

$$\text{age}(x_1) \quad \text{height}(x_2) \quad \text{weight}(x_3)$$

$$= \begin{bmatrix} 4 & 2104+1416+1534+852 & 13 \\ 2104+1416+1534+852 & (2104 \cdot 5) + (1416 \cdot 3) + (1534 \cdot 3) + (852 \cdot 2) & 6 \\ 5 & (2104 \cdot 1) + (1416 \cdot 2) + (1534 \cdot 2) + (852 \cdot 1) & 178 \end{bmatrix}$$

So,

$\theta = (X^T X)^{-1} X^T y$  is the optimal value of  $\theta$  for which we get  $\min_{\theta} J(\theta)$ .

⑥ In training examples,  $m$  features.

Gradient Descent

Normal Equation

- i) Need to choose  $\alpha$
- ii) Needs many iterations.
- iii) Works well even when  $m$  is large.  
[ $O(kn^2)$  time]

$n = 10^6$ , use G.D.

large -

[inverse with  $m=1000$  is still ok, but,  $m=1000$ , it is difficult]  
in that case use G.D.

⇒ Normal Equation for multivariate Linear Regression.

Normal Eqn:

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

- What if  $X^T X$  non-invertible? (singular/degenerate)

( $P^{-1} \theta$  is pseudo inverse).

⇒ If  $X^T X$  is non-invertible, there usually two most common causes for that.

v) case 1) If somehow, in your learning problem, you have redundant prob features

Concrete, if you are trying to predict housing prices and if  $x_1$  is the size of house (in feet) and  $x_2$  is the size in meter, then  $x_1$  and  $x_2$  are linearly dependent ( $x_1 = (3.28)x_2$ ).

In linear algebra, one can see that if two features are related by some linear equation, then the matrix  $(X^T X)$  is non-invertible.

Cause ii) If you are trying to run the learning algorithm with a lot of features.

Concretely,  $m \leq n$ .

Let  $m = 10$ ,  $n = 100$ .  
Now for so,  $\theta$  error. i.e., you are trying to fit 10 training parameters from just 10 training examples.

It is not a good idea.  
In this case, either delete some features

or use regularization!

## ⑩ Classification

### » Logistic Regression

Some examples of classification:

→ Email : spam / not spam?

→ Online Transaction : fraudulent / not fraudulent

→ Tumors : Malignant / Benign?

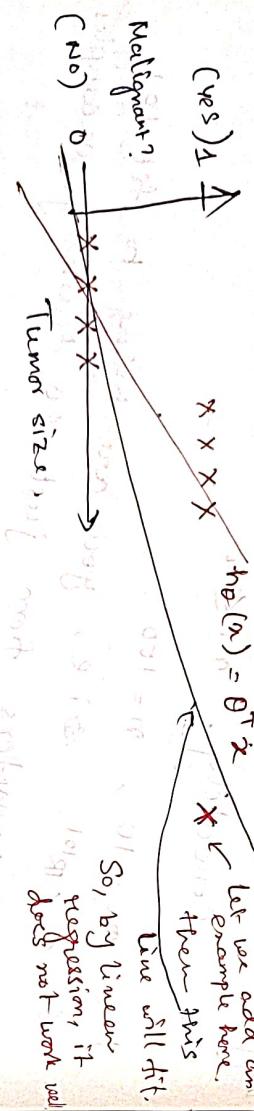
$$y = \{0, 1\}$$

0: "negative class" (e.g., benign tumor)  
1: "positive class" (e.g., malignant)

Let us start for 2 - classes only (0 & 1).

Later we'll discuss multi-class problem.

How do we develop a classification algorithm?



Even though the term "regression" is present with "logistic regression", still, logistic regression is a classification problem.

Classification:  $y = 0$  or  $1$

$$\theta_0(x) \text{ can be } > 1 \text{ or } < 0 \text{ (if use regression)}$$

Logistic Regression has the following property,

$$0 \leq \theta_0(x) \leq 1$$

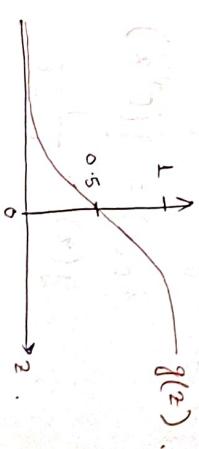
In linear regression,  $\theta_0(x) = \theta^T x$ , where  $\theta_0(x)$  as  $0.5$ .

In logistic " ",  $\theta_0(x) = g(\theta^T x)$ , where  $g(z) = \frac{1}{1+e^{-z}}$  → this is called the sigmoid function or logistic function.

→ You can decide a threshold classifier output  $\theta_0(x)$  as  $0.5$ .  
 if  $\theta_0(x) > 0.5$ , predict "y = 1"  
 if  $\theta_0(x) < 0.5$ , predict "y = 0".

→ It may look like linear regression, in this example, even though it is a classification task we are interested in.

value of  $g(z)$  is also between 0 and 1.



## > Interpretation of Hypothesis Output:

$\rho_{\theta}(x)$  = estimated probability that  $y=1$  on input  $x$

Example:

$$\text{If } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$$

$$\rho_{\theta}(x) = 0.7$$

Tell patient that 70% chance of tumor being malignant.

So, interpret as,

$$\rho_{\theta}(x) = P(y=1|x; \theta) \quad \text{or simply the above probability that } y=1, \text{ given } x, \text{ parameterized by } \theta.$$

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

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

$$\Rightarrow P(y=0|x; \theta) = 1 - \rho_{\theta}(x)$$

## ⑩ Decision Boundary :-

How the hypothesis function looks like particularly when we have more than one feature.

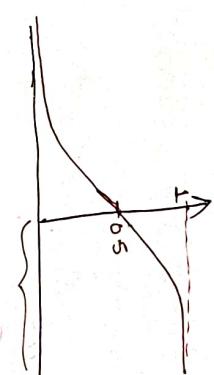
In Logistic Regression,

$$\rho_{\theta}(x) = g(\theta^T x) = P(y=1|x; \theta)$$

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

Suppose, predict "y=1" if  $\rho_{\theta}(x) > 0.5$  and "y=0" if  $\rho_{\theta}(x) \leq 0.5$ .

$$\text{predict "y=0" if } \rho_{\theta}(x) < 0.5 \Rightarrow \theta^T x < 0$$



$$g(x) > 0.5 \text{ when } \theta^T x > 0$$

$$\text{as, } \rho_{\theta}(x) = g(\theta^T x) > 0.5$$

### Decision Boundary:-

Let, we have a training set as shown in the figure and,

$$\text{The straight line is called decision boundary.}$$

$$\text{Let, values of } \theta_0 = -3, \theta_1 = 1 \text{ and } \theta_2 = 1.$$

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

is predicted & when  $y=1$

So, let us see when  $y=0$  is predicted & when  $y=1$  is predicted.

predict, " $y=1$ " if  $-3 + x_1 + x_2 \geq 0$ .  $\left[ \because \theta^T x \geq 0 \text{ for } y=1 \right]$ .

$\Rightarrow x_1 + x_2 \geq 3$ .  $\Rightarrow$  defines a straight line.

So, Decision Boundary is a property of the hypothesis,  $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$  (In this example, not a property of the dataset)

Classification Problem  
So,  $h_{\theta}(x) \geq 0$  when  $x_1^2 + x_2^2 \geq 1$ .

Decision boundary

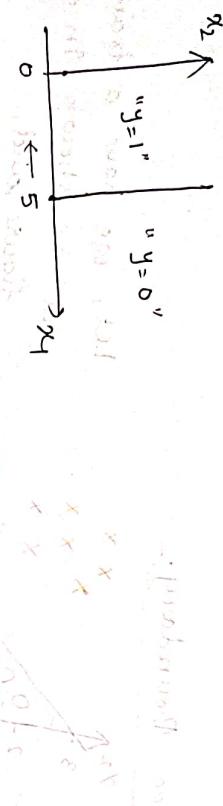
Example:

$$\theta_0 = 5, \theta_1 = -1, \theta_2 = 0.$$

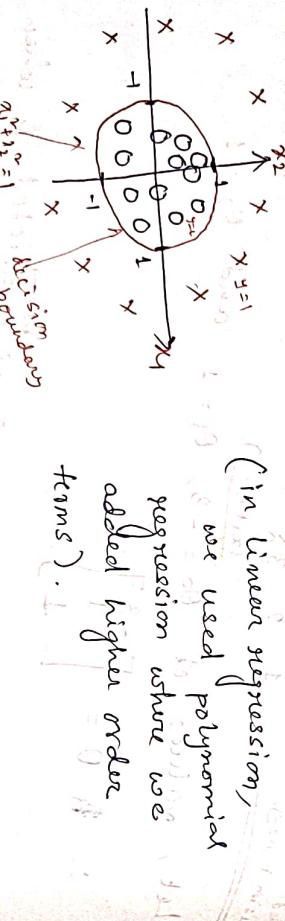
$$h_{\theta}(x) = g(5 - x_1).$$

So,  $h_{\theta}(x) \geq 0$  when  $5 - x_1 \geq 0$

$$\Rightarrow x_1 \leq 5.$$



Non-linear decision boundaries  $\Rightarrow$



Similarly, here also, we use 5 parameters.

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

$$\text{Let, } \theta_0 = -1, \theta_1 = 0, \theta_2 = 0, \theta_3 = 1, \theta_4 = 1$$

After simplifying, we get,  
 $\text{cost } (h_{\theta}(x) - y) = \frac{1}{2} (h_{\theta}(x) - y)^2$

This is the cost I want my learning algorithm to pay if it outputs that value if its prediction is  $h_{\theta}(x)$  and actual label was  $y$ .

Predict "y=1" if  $-1 + x_1^2 + x_2^2 > 0$   
 $\Rightarrow x_1^2 + x_2^2 > 1$ .

So, decision boundary looks like a circle.

Cost functions used of Logistic Regression:-

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

$$\text{in examples: } x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}, x_0 = 1, y \in \{0, 1\}$$

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

How to choose parameters  $\theta$ ?

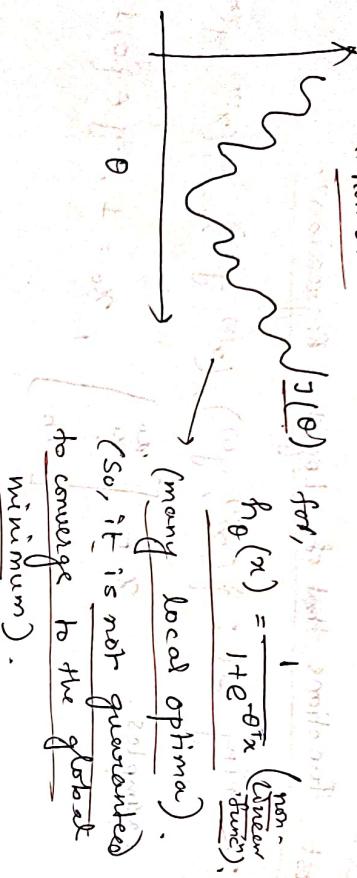
Cost function of Linear Regression

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

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

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

For the Logistic Regression, if we would use the same cost function, we get non-convex function of the parameter's data. "non-convex"



'Convex' is bow-shaped where it converges to global minimum, where is only one global minimum).

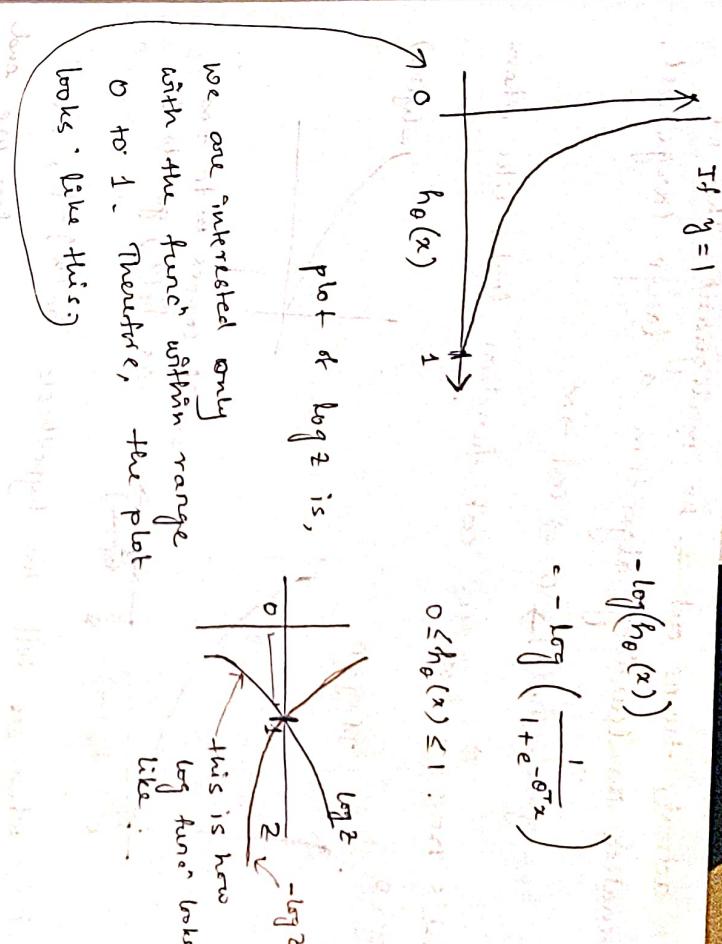
for the non-linear function  $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$ ,  $J(\theta)$  ends up being a non-convex function if you were to define it as a square cost function.

## ② Logistic Regression Cost Function:-

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

The cost or penalty that the algorithm pays, if it upwards the value of  $h(x)$ .

Let's plot the function:



So, the hypothesis  $h_\theta(x)$  with  $y=1$  is basically

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

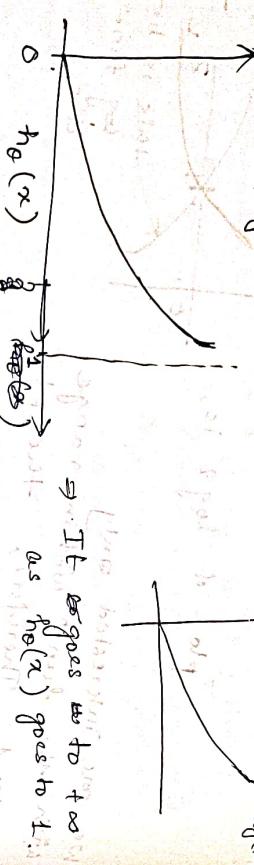
Now, when we write  $h_\theta(x)=0$ , this means,  $P(y=1|x; \theta)=0$  (when  $y=1$ ). That is, we are saying that, say for a medical patient, probability of being a malignant patient is 0. This depicts as being a malignant patient i.e., we are making it sure that level of certainty, i.e., we are making the patient the patient is not malignant. But, let the patient

in actual, turns out to be malignant. Then we penalize the learning algorithm a large large cost.

so,  $y=1, h_\theta(x) \rightarrow 0 \Rightarrow \text{cost} \rightarrow \infty$

Let's now see how the cost function looks like when  $y=0$ .

If  $y=0$



i.e., when we tell the hypothesis  $h_\theta(x)=1$ , i.e.  $P(y=0|x; \theta)=1$  with certainty, then also we end up with a large cost.

Similarly, when  $y=0$  and  $h_\theta(x)=0$   $\Rightarrow$  cost = 0

### Cost Function of Logistic Regression

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

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

for a "classification" problem,  $y$  always  $\in \{0, 1\}$ .

Concretely, we can write,

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

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

This cost function is derived from statistics using the principle of maximum likelihood estimation, which is an idea in statistics of how to efficiently find parameters' data for different models. Also, it has a nice property that it is convex.

To fit parameters  $\theta$ :

$$\min_{\theta} J(\theta) \Rightarrow \text{Get } \theta$$

To make a prediction given new  $x$ :

$$\text{Output } h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} \text{ with the values of }$$

So, let's see how to minimize  $J(\theta)$ .

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

Want  $\min_{\theta} J(\theta)$ :

$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \\ \text{(Simultaneously update all } \theta_j) \end{array} \right.$$

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

If you have  $n$  features, then

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\begin{bmatrix} \text{here,} \\ h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} \end{bmatrix}$$

» feature scaling makes gradient descent run faster.

## ① Advanced optimization Idea

### » Optimization Algorithm

cost function  $J(\theta)$ . Want  $\min_{\theta} J(\theta)$  with no local minima.

Given  $\theta$ , we have code that can compute

- $J(\theta)$
- $\frac{\partial}{\partial \theta_j} J(\theta)$  (for  $j=0, 1, \dots, n$ )

1. Gradient descent: same idea as in linear regression, but along  $\theta$ .

Repeat {

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

}

Advantages:

- No need to manually pick  $\alpha$ .
- often faster than gradient desc.

2. Conjugate gradient

- Disadvantage:
  - More complex.

3. BFGS

4. L-BFGS

① Example : ( how to use these algorithms )

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

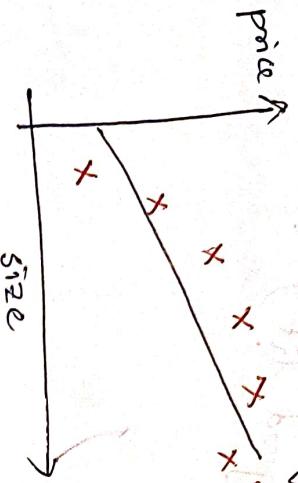
$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

<use similar code to use an optimized algorithm>

## ① The problem of overfitting :-

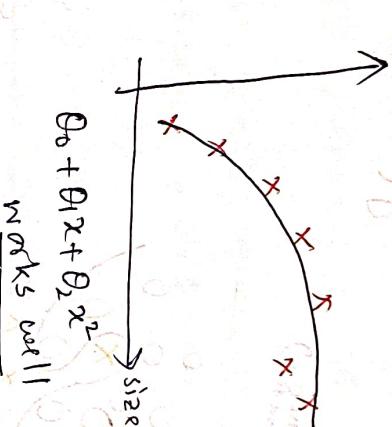
Example: Linear Regression (housing prices)



(not a very good model)

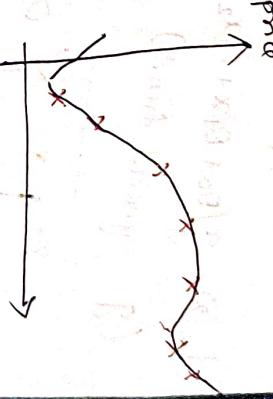
"High bias"

does not fit ~~properly~~.



works well

seems to have  
good job, but  
actually not.



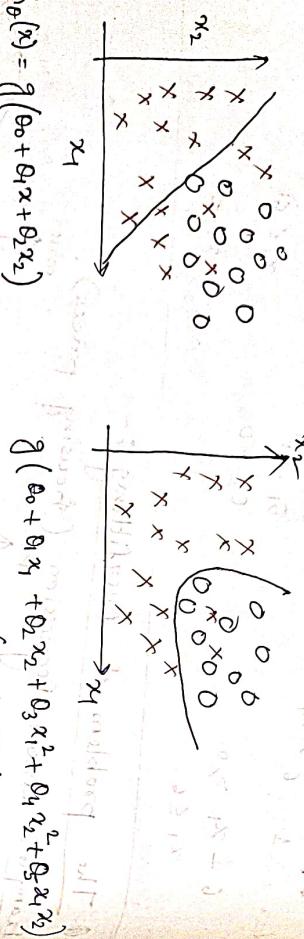
overfitting  
"high variance".

The handwritten note below provides intuition for this phenomenon:

"The intuition is that, if we're fitting such a high order polynomial, then the hypothesis space of possible functions is very large, and we do not have enough data to constrain it to give us a good hypothesis. So, it almost fit any func" and this "face of overfitting."

» The problem of overfitting comes when we have too many features, the learned hypothesis may fit the training set very well ( $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i(x^{(i)}) - y^{(i)})^2$ ) but fail to generalize to new examples (predict prices on new examples).

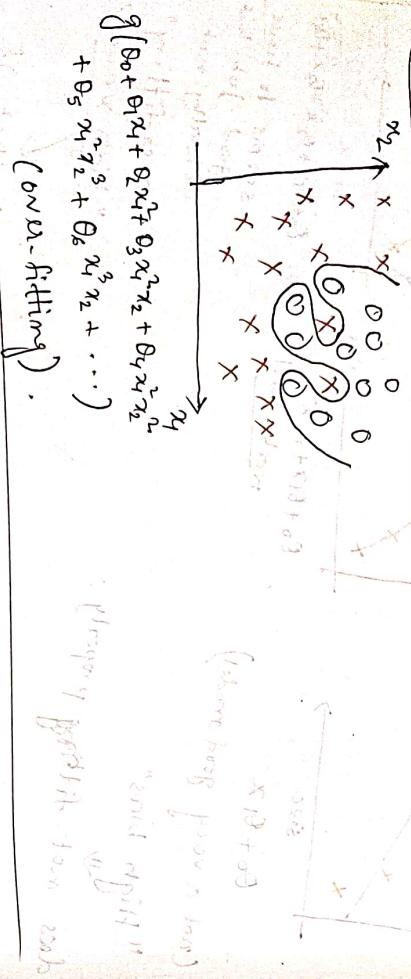
» Some thing can be applied for logistic regression as well.



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

$$(g = \text{sigmoid func.})$$

(good fit)



$$J(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots)$$

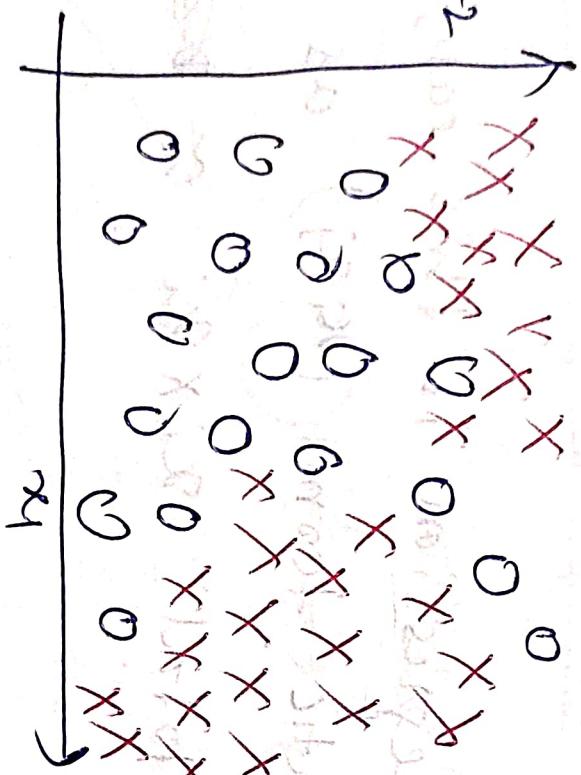
(more terms)

more terms

## Non-linear Hypothesis

### Neural Networks

Consider a Non-linear classification:-



Let, you have a training set like this. If you want to apply logistic regression, one thing you could do is,

apply logistic regression with lot of non-linear features like  $g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \dots)$

$$\theta_0 x_1^0 x_2^0 + \theta_1 x_1^1 x_2^0 + \theta_2 x_1^0 x_2^1 + \dots$$

But, this particular method works only when you have two features  $x_1$  &  $x_2$  and all the polynomial terms of  $x_1$  &  $x_2$  in  $g$  (which is the sigmoid function).

But, we may have lot more features than just two.

Let, in a housing problem,

$$x_4 = \text{size}$$

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

$$x_3 = \# \text{floors}$$

$$x_4 = \text{age}$$

$$x_{100}$$

Let, we are taking an example of identifying a car from an image.

If we chose only two pixel areas of  $50 \times 50$  pixel size of each region, then we have  $50 \times 50 = 2500$  pixels. I.e.,  $n = 2500$  (7500 if RGB).

$$\mathbf{x} = \begin{bmatrix} \text{pixel 1 intensity} \\ \text{pixel 2} \\ \vdots \\ \text{pixel 2500} \end{bmatrix}$$

Now, if we use Logistic regression in this case, including all the quadratic terms ( $x_i x_j$ ) as features, then there are total  $\frac{n(n+1)}{2} \approx 300000$  features. million features.

Origins of Neural Network :- Algorithms that try to mimic the brain.

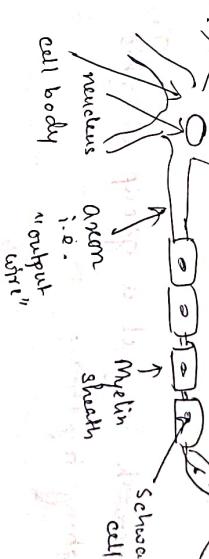
→ Was very widely used in 80s and early 90s; popularity diminished in late 90s.

Recent Resurgence :- State-of-the-art techniques for many applications as computers become faster. NN has complex algorithm.

Neuron in brain "input wires"

"output wires"

"output wire"



Neurons communicate with each other by little pulse of electricity.

⇒ Neuron Model : Logistic unit

$$\text{"output wires"} \rightarrow h_0(\mathbf{x})$$

$$\begin{array}{l} \text{"input wires"} \\ \text{---} \\ (\mathbf{x}_1) \quad (\mathbf{x}_2) \end{array} \rightarrow \text{bias unit}$$

Sometimes,  $x_0$  is also drawn.  $x_0$  is called the "bias unit". because,  $x_0$  is already equal to 1.

⇒ "Activation func" is a terminology used in Neural Network which is a term or function of non-linearity

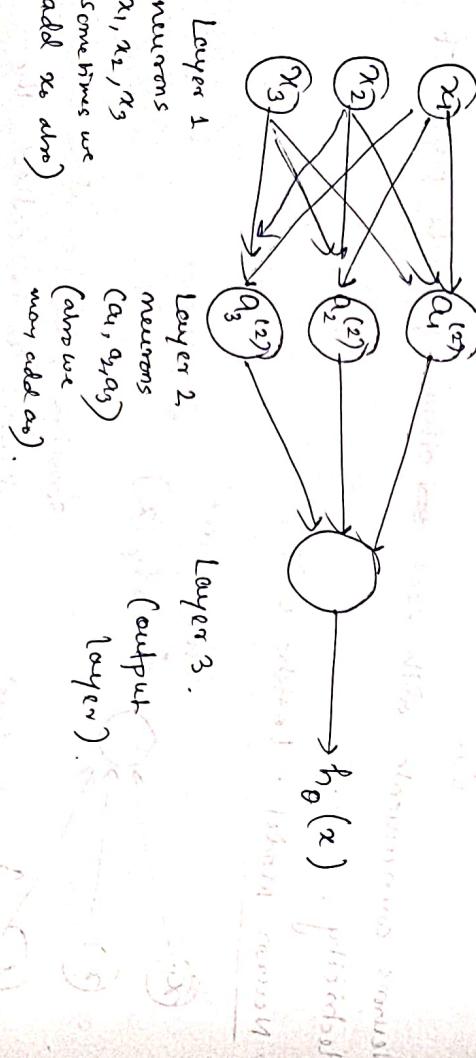
$$g(z) = \frac{1}{1+e^{-z}}$$

⇒ "Weight" of a model is same thing as the parameters ( $\theta_0, \theta_1, \theta_2, \dots$ ) of a model.

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

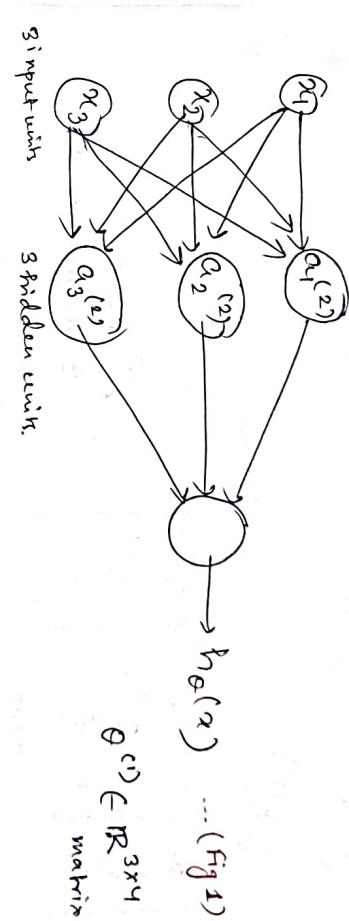
So, the above diagram (prev. page) represent a single neuron.

⇒ What a neural network is, just a group of this different neurons strong together.



⇒ a hidden layer of values you don't get to observe in the training set-up.

On the above example, we have only one hidden layer. But there can have more than one hidden layer.



$\alpha_i^{(j)}$  = "activation" of unit  $i$  in layer  $j$

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

$$\alpha_1^{(2)} = g(\theta_{10}x_0 + \theta_{11}x_1 + \theta_{12}x_2 + \theta_{13}x_3) \quad [\text{sigmoid or activation function}]$$

$$\alpha_2^{(2)} = g(\theta_{20}x_0 + \theta_{21}x_1 + \theta_{22}x_2 + \theta_{23}x_3)$$

$$\alpha_3^{(2)} = g(\theta_{30}x_0 + \theta_{31}x_1 + \theta_{32}x_2 + \theta_{33}x_3)$$

If "also logistic function".

Layer 1 is input layer  
Layer 2 is hidden layer

⇒ If network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j+1$ , then  $\boldsymbol{\theta}^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .

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

~~=  $\theta_{10}^{(1)}$~~

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

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

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

$$\begin{aligned} x &= \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}, \quad z^{(2)} \text{ is a } 3\text{-dimensional vector.} \\ &\uparrow \qquad \qquad \qquad \uparrow \\ \text{feature vector.} & \text{vector of } z\text{-values.} \end{aligned}$$

$$z^{(2)} = \theta^{(1)} x + \text{bias} = \theta^{(1)} x + b$$

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

$\in \mathbb{R}^3$

(i.e., a 3-dimensional vector).

Here  $g$  is the activation function of the elements of  $a^{(2)}$ .

$z^{(2)}$

we can consider,

$$a^{(1)} = x \cdot \text{bias} + \theta^{(0)} x + \theta^{(1)} x + \theta^{(2)} x + \dots + \theta^{(n)} x$$

so,  $a^{(1)}$  is activation's input layer.

There is a 'biased unit' in every layer.

$$z^{(3)} = \theta^{(2)} a^{(2)}. \quad f_\theta(x) = a^{(3)} = g(z^{(3)})$$

» We start off the propagation from input units and then we sort of forward-propagate that to the hidden layer and compute the activations at the hidden layer and forward propagate that and come find output layer.  
So, Input layer  $\rightarrow$  hidden layer  $\rightarrow$  output layer.  
So, we call it Forward Propagation.

The hidden layer or neural network computes logistic regression and get the hypothesis function

$$f_\theta(x) = g(\theta^{(2)} x)$$

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

The only difference from logistic regression is that, in logistic regression, we use  $x_1, x_2, x_3$  directly instead of  $a_1^{(1)}, a_2^{(1)}, a_3^{(1)}$ . That means layer 2 is learning by their own from a layer 1 and computes  $h_\theta(x)$ .

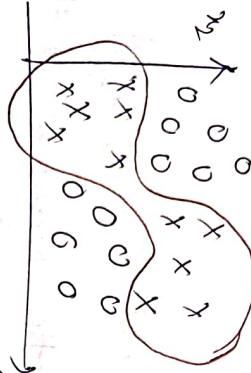
Network architecture means how the net neurons are connected to each other.

## >> Example 1 :-

Non-linear classification example : XOR/ XNOR .  
 $x_1, x_2$  are binary (0 or 1).



Let us have positive bunch of positive examples (denoted by 'x') and negative examples (denoted by 'o')



So, our target is to learn a non-linear decision boundary that separates the positive and the negative examples.

$$\begin{array}{c} y \\ \text{x} \\ \text{o} \end{array}$$

$$\begin{array}{c} x_1 \\ x_2 \end{array}$$

$$\begin{array}{c} y \\ \text{x} \\ \text{o} \end{array}$$

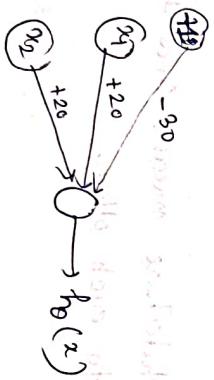
$$\begin{array}{c} x_1 \\ x_2 \end{array}$$

Simple example: AND

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

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

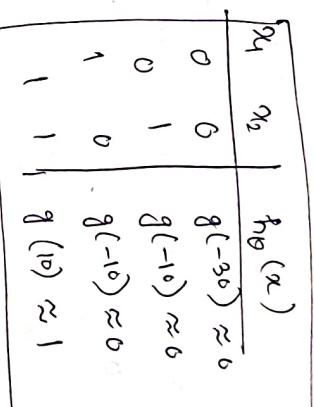
$$= \text{NOT}(x_1 \text{ XOR } x_2)$$



Can we get a one unit network to compute this logical AND operation? In order to do so, I am going to draw the bias unit. (i.e.,  $\oplus 1$  unit).

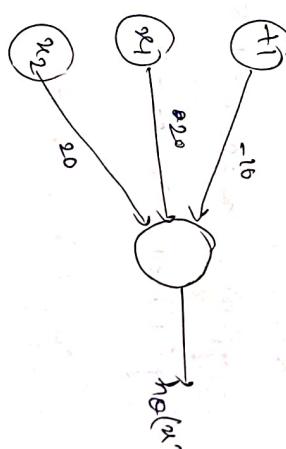
Let the parameters are,  $-30$ ,  $+20$ ,  $+20$  resp.  
 $\therefore h_0(x) = g(-30 + 20x_1 + 20x_2)$

$$\begin{matrix} x_1 & x_2 & h_0(x) \\ 0 & 0 & g(-30) \approx 0 \\ 0 & 1 & g(-10) \approx 0 \\ 1 & 0 & g(10) \approx 1 \\ 1 & 1 & g(30) \approx 1 \end{matrix}$$



i.e.  $h_0(x) \approx x_1 \text{ AND } x_2$  that our neural network computes.

Ex 2 :-



$$\begin{matrix} x_1 & x_2 & h_0(x) \\ 0 & 0 & g(-10) \approx 0 \\ 0 & 1 & g(10) \approx 1 \\ 1 & 0 & g(10) \approx 1 \\ 1 & 1 & g(30) \approx 1 \end{matrix}$$

So, it is  $x_1 \text{ OR } x_2$ .

(Ans)

$$h_0(x) = g(-10 + 20x_1 + 20x_2)$$

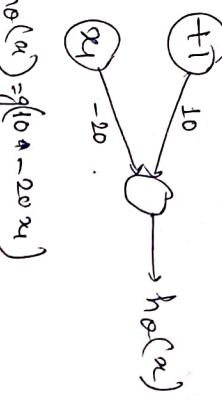
$$= g(-10) \approx 0, \text{ when } x_1, x_2 = 0$$

$$= 1 \cdot g(10), \text{ when } x_1 = 0 \text{ & } x_2 = 1$$

$$= g(10), \quad " \quad x_1 = 1 \quad x_2 = 0$$

$$= g(30) \quad " \quad x_1 = x_2 = 1$$

### Ex 3: Negation



Ex 4:  $(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$ .

Logical function is,

= 1 if and only if  $x_1 = x_2 = 0$

<u>AND</u>	
$x_1$	$\sim x_1$
$x_2$	$\sim x_2$
0	0
1	0
0	1
0	1

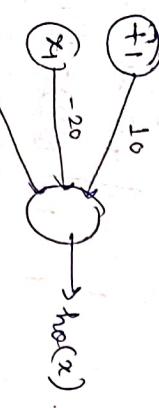
$$\boxed{h_0(x) = g(10 + -20x_1 + -20x_2)}$$

$$= g(-30) \approx 0 \Rightarrow x_1 = x_2 = 1$$

$$= g(-10), x_1 = 0, x_2 = 0 \approx 0$$

$$= g(-10), x_1 = 0, x_2 = 1 \approx 0$$

$$= g(10) \approx 1, x_1 = x_2 = 0$$



i.e.  $x_1 \text{ AND } x_2$  : Putting  $x_1$  AND  $x_2$ ,  $(\text{NOT } x_1)$  AND  $(\text{NOT } x_2)$

and  $x_1$  OR  $x_2$  together.



<u>AND OR</u>	
$x_1$	$\sim x_1$
$x_2$	$\sim x_2$
0	0
1	0
0	1
1	1

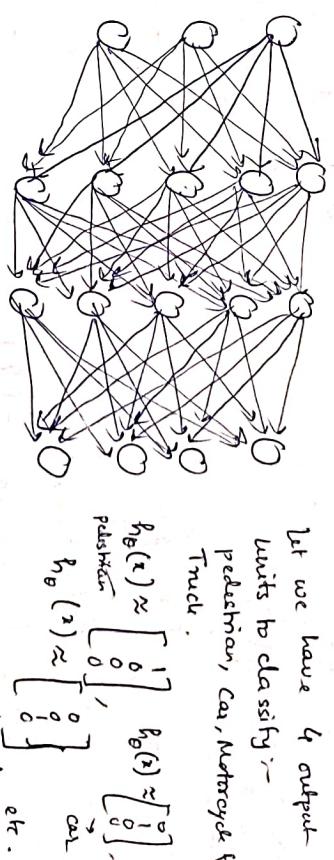
Ans:  $a_1^{(1)} : \text{AND}, a_2^{(1)} : (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

1.  $x_1 \text{ XOR } x_2$
- $= (x_1 \text{ AND } x_2) \text{ OR } [(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)]$

$\Downarrow$   
 $a_1^{(3)} \Rightarrow \text{output layer}$

$\Downarrow$   
 $a_2^{(3)} \Rightarrow \text{hidden layer}$

### Multiclass classification in Neural Network



Let we have 4 output units to classify :-

Pedestrian, Car, Motorcycle & Truck.

$h_0(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ ,  $h_0(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$

$h_0(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ ,  $h_0(x) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

etc. etc.

Hand written digits : multiclass classification

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

$y^{(i)}$  is one of  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$  &  $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \in \mathbb{R}^4$ .

pedestrian, car, Motorcycle, Truck  
(4 dimensional vector).



$$\theta_{j_1}^{(1)} \Rightarrow \theta_{j_0}^{(2)} x_0 + \theta_{j_1}^{(2)} x_1 + \dots$$

[ $\theta_{j_0}$  bias unit,  $x_0$  feature +  $x_1$  bias unit]

### ④ Back Propagation Algorithm :

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

The ' $\delta$ ' term is, in some sense, going to capture the error in activation function of that node, so how we might wish the activation of that node was slightly different.

For each output ~~feature~~ unit (Layer  $L=4$ )

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

[ $y_j$  is the actual value observed in our training example].

So,  $a_j^{(4)}$  can also be written as,

$$(h_\theta(x))_j$$

So,  $\delta$  is the difference b/w hypothesis & actual value.

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

Now, each  $\delta^{(4)}$ ,  $a^{(4)}$  &  $y$  is a vector whose dimension is equal to the number of output units in our network.

$$\delta^{(3)} = (\theta^{(3)})^\top \cdot \delta^{(4)} \cdot g'(z^{(3)}).$$

elementwise multiplication operation.

$$\delta^{(2)} = (\theta^{(2)})^\top \cdot \delta^{(3)} \cdot g'(z^{(2)}).$$

Here,  $g'(z^{(3)})$  formally is actually the derivative of the activation function  $g$  evaluated at the input values given by  $z^{(3)}$ .

After calculating using calculus, we get,

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

$$= a^{(3)} \cdot (1 - a^{(3)})$$

$$g'(z^{(2)}) = a^{(2)} \cdot (1 - a^{(2)})$$

$a$  is the vector of activations

1 is the vector of 1s.

⇒ So there is no  $\delta^{(1)}$  term, as the 1st layer is the input layer and the features we observed in our training set don't have any error associated with it.

⇒ So, we start computing the ' $\delta$ '-term from the output layer, and go back to the previous layers. So, we call it as back propagation.

→ Though mathematically complicated, but the following is possible to prove.

$$\frac{\partial}{\partial \theta_{ij}^{(n)}} J(\theta) = a_j^{(n)} \delta_i^{(n+1)} \quad [\text{ignoring regularization term } \lambda; \text{ or if } \lambda=0]$$

[Next continued after Decision page 34 pages].

## Bias, Variance, Overfitting & Underfitting

~~Distance between predicted value & actual value, i.e., "how far we have predicted the values from actual values.~~

High Bias: If the average prediction value is far away from actual values.

Low Bias: Distance between predicted & actual values are minimal.

Let,  $y = f(x) + e$ ,  $e$  is error.

Expected squared error at point  $x$  is  $E[(y - \hat{f}(x))^2]$

$$E[\text{Err}(x)] = E[(y - \hat{f}(x))^2]$$

$$= E[(\hat{y} - f(x)) + (f(x) - \hat{f}(x))^2]$$

$\hat{f}(x)$  is predicted value.  
 $y$  is original value.

$$\begin{aligned}
 &= E \left\{ (y - f(x))^2 + (f(x) - \hat{f}(x))^2 + 2[y - f(x)][f(x) - \hat{f}(x)] \right\} \\
 &= E[(y - f(x))^2] + E[(f(x) - \hat{f}(x))^2] + 2[E[y] - E[f(x)]] \\
 &\quad - [E[y] + E[f(x)]]^2 + [E[(y - f(x))(\hat{f}(x) - f(x))]] \\
 &\quad - E[\hat{f}^2(x)] + E[f(x)\hat{f}(x)].
 \end{aligned}$$

$f(x)$ : true func

$\epsilon$ : Error

$\hat{y}$ : expected value from dataset D

$\hat{y}$ : predicted value of target variable

$h_D$ : predicted value learned from many datasets D

(\*) for  $\hat{y}$  not following

true func been

with respect to noise

learned from D the best fit function

long term expectation

learner's prediction on this x averaged over many datasets D.

(\*) can be written as,

$$\checkmark E_{D,\epsilon} \{(f - \hat{y})^2\}$$

$$= E \{ ([f - h] + [h - \hat{y}])^2 \}$$

$$= E \{ [f - h]^2 + [h - \hat{y}]^2 + 2 [f - h] [h - \hat{y}] \}$$

$$= E \{ [f - h]^2 + [h - \hat{y}]^2 + 2 [fh - f\hat{y} - h\hat{y} + h\hat{y}] \}$$

$$= E \{ (f - h)^2 \} + E \{ (h - \hat{y})^2 \} + 2(E[f\hat{y}] - E[f\hat{y}]) -$$

$$E[\hat{y}\hat{y}] + E[h\hat{y}]$$

$$= E[(f-h)^2] + E[(h-\hat{y})^2]$$

Squared difference between best possible prediction for  $x$ , i.e.,  $f(x)$  and our long-term expectation for what the learner will do if we averaged over many datasets, i.e.,  $E[h(x)]$ .

Then,

$$\text{Variance} = E[(\hat{y} - E[\hat{y}])^2]$$

The expected value of the square of the difference between predicted values and the expected value of the predicted values.

### Continued Neural Network

#### Backpropagation Algorithm

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

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

This is used to compute

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

So, Error can be further ~~presented~~ decomposed to,

$$E_{\text{err}}(\theta) = \underbrace{\text{Bias}^2}_{\text{Bias error}} + \underbrace{\text{Variance}}_{\text{Variance error}} + \underbrace{\text{Irreducible Error}}_{\text{Irreducible error}}$$

① Variance is the variability of the model that how much it is sensitive to another subset of the training dataset.

If  $y$  is actual value.  
 $\hat{y}$  is predicted value of target variable.

reducible error.

These  $\Delta$  deltas (capital) are going to be used accumulators that will slowly add things to in order to compute these partial derivatives.

Set  $a_{ij}^{(0)} = 0$  (for all  $i, j$ )

$a^{(0)} = x^{(0)}$

For  $i = 1$  to  $m$   
Perform forward propagation to compute

$a^{(l)}$  for  $l = 2, 3, \dots, L$  is taking it by

Using  $y^{(i)}$ , compute

$$\delta^{(L)} (= a^{(L)} - y^{(i)})$$

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

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l-1)} + a_j^{(l)} \delta_i^{(l+1)}$$

$$\Delta_{ij}^{(0)} := \frac{1}{m} \Delta_{ij}^{(1)} + \lambda \theta_{ij}^{(0)}, \text{ if } j \neq 0$$

(when we do not have regularization)

$$\frac{\partial J(\theta)}{\partial \theta_{ij}} = \Delta_{ij}^{(1)} \quad [\text{the proof is pretty complicated}]$$

This is back propagation algorithm.

## Clustering

In clustering, we would like to have an algorithm that automatically group the data into coherent clusters (i.e., coherent clusters for us).

K-means clustering is, by far, the most popular clustering algorithm.

### Partitional clustering

Or is a type of clustering that divides the dataset into non-overlapping subsets (or partitions) such that each data point is included in exactly one subset.

example: K-means clustering, k-medoids, fuzzy c-means etc.

### K-means clustering :- An algorithm to achieve

Partitional clustering. It aims to partition the data into  $k$ -clusters in which each data point belongs to the cluster with the nearest mean (center of the cluster), serving as a prototype of the cluster.

1. Initialise centroid in each cluster
2. Calculate Euclidean distance of a point to the centroid,
3. Find min of the dist calculated value, and
4. Update the centroid by taking mean
5. Follow step 2 to 4.

Example:

sl. no.	age	amount
C1	20	500
C2	40	1000
C3	30	800
C4	18	300
C5	28	1200
C6	35	1400
C7	45	1800

Step 1: Let  $K=2$ . We'll create two clusters centroids randomly.  
 Step 2: Initialize centroids with random values. For example,  $c_1 = (35, 900)$  and  $c_2 = (20, 500)$ .

Step 3: update centroid of  $K=2$  cluster. Centroid of  $K=2$  cluster is,  $\left(\frac{40+30}{2}\right), \left(\frac{1000+800}{2}\right) = (35, 900)$

$c_3$  is in  $K=2$  cluster.

Step 4: calculate distance between  $c_3$  and  $c_1$ ,  $c_2$  and put in the corresponding cluster for which we get lesser value.

for  $c_4$  in  $K=2$

for  $c_4$  in  $K=1$

Step 5: calculate distance between  $c_4$  and  $c_1$ ,  $c_2$  and put in the corresponding cluster for which we get lesser value.

continue the process.

>> Hierarchical clustering

Instead of dividing the dataset into a pre-determined number of clusters at once, hierarchical clustering builds up a hierarchy of clusters either

through a divisive method (top-down approach) or an agglomerative method (bottom-up approach). It helps in making relation between two clusters.

For instance,  $c_3$  is  $k_2$ , so it must be included in both clusters. Hence,  $d = \sqrt{(30-40)^2 + (800-1000)^2} = 200.24$ .

As,  $200.24 < 300.00$ ,  $c_3$  belongs to both of the clusters.

## Agglomerative clustering

- Initially consider every data point as an individual cluster and at every step, merge the nearest pairs of the cluster. (bottom-up).

- At first, every dataset.

### steps:

1. calculate the similarity of one cluster with all the other clusters (calculate proximity matrix or other distance metric).

2. consider every data point as an individual cluster.

3. Merge the clusters which are slightly similar or close to each other.

4. Recalculate (update) the proximity matrix

for each cluster.

5. Repeat the steps 3 & 4 until only a single cluster remains.

Example: Consider 5 data points, 1, 2, 3, 4, 5.

Let's make a distance matrix between pair of units that you want to cluster.

A distance matrix will be symmetric (because, the distance between  $x$  and  $y$  is the same as the distance from  $y$  to  $x$ ). and will have 2 zeros on the diagonal (because, every item is distance zero from itself).

cluster	1	2	3	4	5
1	0	35	0	0	0
2	35	0	9	0	0
3	0	9	0	7	0
4	0	7	0	0	5
5	0	0	5	11	0
	35	11	9	7	5

minimum is 2 i.e., merging of 3 & 5. new cluster is '35'. Remove '3' & '5' & replace it by '35'.

Since we are using complete linkage clustering, the distance between "35" and every other item

is the maximum of the distance between this (item & 3) and (item & 5). For example,

(item & 2) distance between

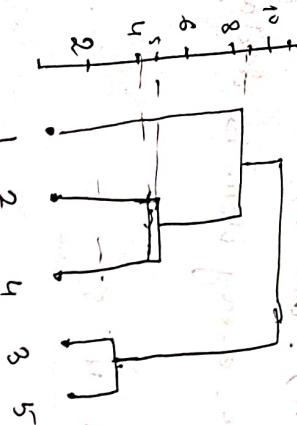
"35" & "1" is

$$\max(d(1,3), d(1,5)) = \max(3, 11) = 11.$$

So, now, we have to update the distance matrix.

cluster	1	2	3	4
1	0	35	0	0
2	35	0	9	0
3	0	9	0	7
4	0	7	0	0

Continuing this way, after 6 steps, everything is clustered.



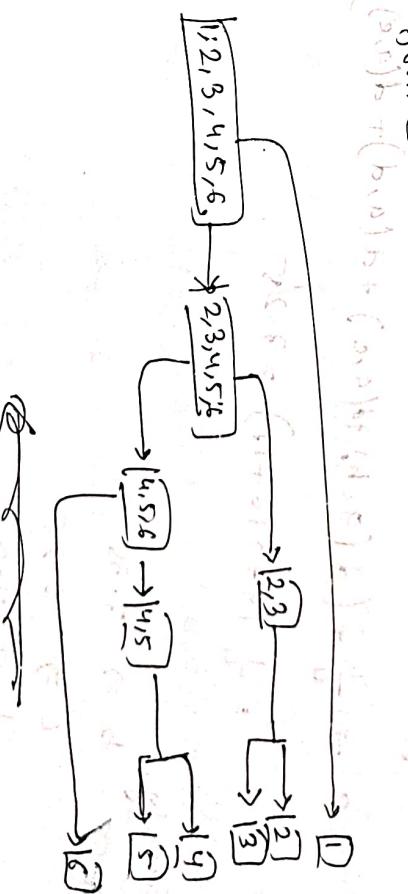
This is called dendrogram  
(a tree-like structure of clustering)

P.E.g.: cluster height of '35' is 2.

Let us now see single linkage dendrogram (below) for the same distance matrix.  
It also starts from '35', but the distance between '35' and each item is now the minimum of  $d(a,3)$  &  $d(a,5)$ .  $d(3,5) = 3.5$

- Divisive clustering?
  - Consider all the data points as a single cluster.
  - Split into clusters using any flat-clustering method, say k-means.
  - Choose the best cluster among the clusters to split further, choose the one that has the largest sum of squared error.

Step 1: Repeat steps 2 & 3 until a single cluster is formed.



top-down approach

Example of divisive clustering

Given dataset with distance matrix:

	a	b	c	d	e
a	0	9	3	6	11
b	9	0	7	5	10
c	3	7	0	9	2
d	6	5	9	0	8
e	11	10	2	8	0

single linkage

Step 3: Initially  $C_a = \{a; b; c; d; e\}$ .

Step 2:  $C_i = C_a$  and  $C_j = \emptyset$ .

Step 3: Initial Iteration

Let us calculate the average dissimilarity of the objects in  $C_a$  with the other objects.

Q. Average dissimilarity of  $a$  is,

$$D_a = \frac{1}{3} (d(a,c) + d(a,d) + d(a,e)) - \frac{1}{3} (d(a,b))$$

$$= \frac{20}{3} - 9 = -2.33.$$

Similarly,

$$D_c = \frac{1}{3} (d(c,a) + d(c,d) + d(c,e)) - \frac{1}{3} (d(c,b))$$

$$= \frac{14}{3} - 7 = 2.33$$

$$D_d = \frac{1}{3} (d(d,a) + d(d,c) + d(d,e)) - \frac{1}{3} (d(d,b))$$

$$= 0.67$$

$$D_e = \frac{1}{3} (d(e,a) + d(e,c) + d(e,d)) - \frac{1}{3} (d(e,b))$$

Here,  $D_a$  is the largest and  $\geq 0$ .

So, move  $a$  to  $C_j$ .

we now have,

Take item of highest dissimilarity. Here  $b & e$ . choose anyone, say,  $b$ , it will be moved.

We move  $b$  from  $C_j$  to  $C_i$ .

We now have,

$$C_i = \{a, c, d, e\} \quad \& \quad C_j = \{b\}$$

Step 4: Remaining iterations.

$$C_i = \{a, c, d, e\}, \quad C_j = \{b\}$$

(i) 2nd iteran :-

$$D_a = \frac{1}{2} (d(a,c) + d(a,d)) - \frac{1}{2} (d(a,b))$$

$$= -0.5$$

$$S_1 = -13.5$$

$$S_2 = \frac{1}{2} (d(c_1, a) + d(c_1, e)) - \frac{1}{2} (d(c_2, b) + d(c_2, d))$$

$$= \frac{13}{2} - \frac{18}{2} = -2.5$$

All are negative so we stop and from the clusters  $c_1$  &  $c_2$ .

$c_1 = \{a, c\}$  and  $c_2 = \{b, d\}$ .

Step 5: To divide  $c_1$  &  $c_2$ , we compare their diameters.

$$\text{diameter}(c) = \max \{d(a, c), d(a, e), d(c, e)\}$$

$$\text{diameter}(c_1) = \max \{3, 11, 2\} = 11$$

$$\text{diam}(c_2) = \max \{d(b, d)\} = 5$$

As,  $11 > 5$ , consider  $c_1$  cluster has the initial cluster & start dividing. (child and repeat the same process).

### Bias:-

Bias is the difference between the average predict of our model and the expected value which we are trying to predict.

Variance :- Variance is the variability in the model prediction, meaning how much the prediction will change if a different training dataset is used.

### K-Nearest Neighbor

1. calculate Distance

(you can use Euclidean distance)

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

2. Select K-nearest neighbor. (you can tune K-value for getting good result).
3. Majority voting.

### Distances

- Euclidean:  $d(p, q) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$
- Manhattan:  $d(p, q) = |x_2 - x_1| + |y_2 - y_1|$
- Minkowski:  $d(p, q) = \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$

If  $p=2$ , it becomes Euclidean.  
If  $p=1$ , it becomes Manhattan.

Example:

Brightness	Saturation	Class
40	20	Red
50	50	Blue
60	90	Blue
70	25	Red
60	70	Blue
60	10	Red
25	80	Blue

Qsn is,

Brightness	Saturation	Class
20	35	?

1. calculate distance (say, Euclidean).

$$d_1 : \sqrt{(40-20)^2 + (20-35)^2} = 25$$

$$d_2 : \sqrt{(50-20)^2 + (50-35)^2} = 33.54$$

$$d_3 : \sqrt{(60-20)^2 + (90-35)^2} = 68.01$$

$$\begin{aligned} d_4 &= 10 \\ d_5 &= 61.03 \\ d_6 &= 47.17 \\ d_7 &= 45. \end{aligned}$$

2. Arrange them in ascending order:

Brightness	Saturation	Class	Distance
10	25	Red	10
40	20	Red	25
50	50	Blue	33.54
25	80	Blue	45
60	10	Red	47.17
70	70	Blue	61.03
60	90	Blue	68.01

Let's choose  $k = 5$ . So, choose the first five rows of the above table.

For 5-nearest neighbor, the majority class is 'Red'.

?	35	(Red)
20	35	Red

→ How to choose  $k$ ?

1. choosing a very low value will most likely lead to inaccurate prediction.

2. The commonly used value of  $k$  is 5.

3. Always use odd number as the value of  $k$  (to avoid ties, prevent bias to specific classes).