

My  
Wing  
know  
our score.

## Week 1

### Machine Learning Definition

- ~~Arthur's~~ old version: field of study that gives computers the ability to learn without being explicitly programmed
- A computer program is said to <sup>new version</sup> learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.

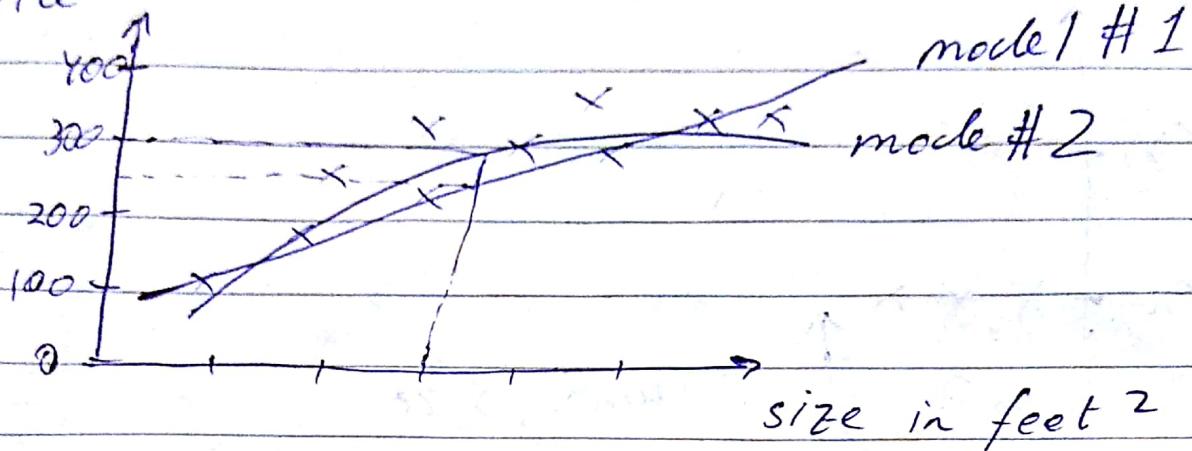
### Machine Learning algorithms (main)

- Supervised learning
- unsupervised learning

Others types of machine learning algos: Reinforcement learning, recommender systems

## Supervised learning ~~example~~ algo example:

Price



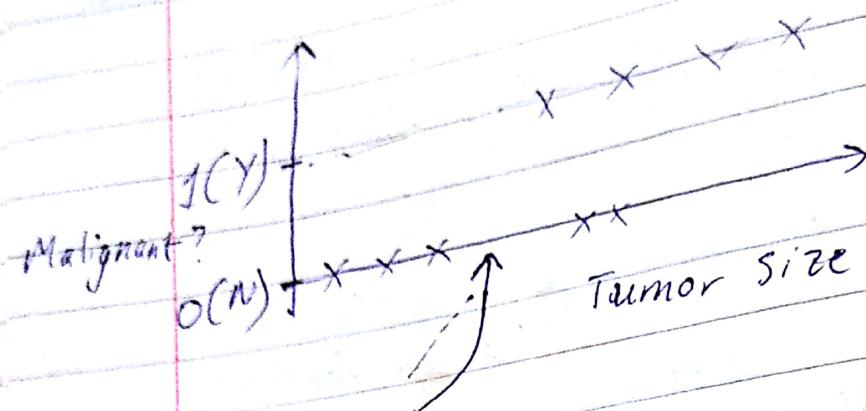
^ Housing price in a particular state  
8 2 learning algorithms are applied

Supervised learning - where the "right answers" to a certain task are present in the data set & algo (using the provided data) tries to approximate the price ~~for~~ for a certain house

Also called a regression problem -  
predict continuous valued output (price in this example)

dangerous  
harmful  
harm

Breast cancer example (malignant, benign)



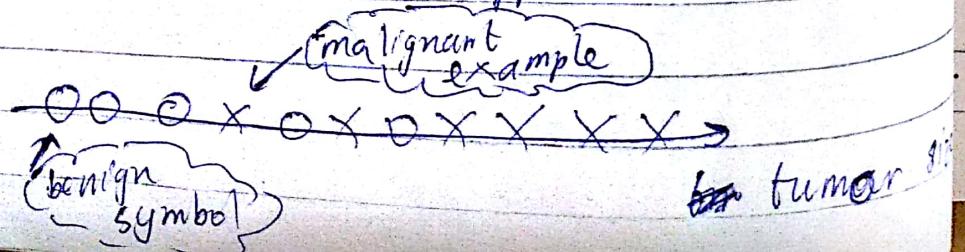
Classification problem example -  
discrete valued outputs (~~can~~ can be more than 2 values )

↳ suppose there are 3 types of cancers : 0, 1, 2, 3

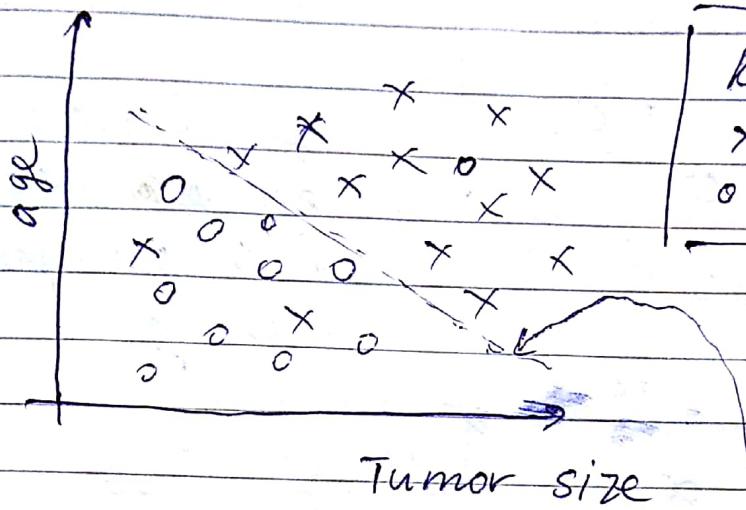
benign      ↑      ↑      ↑      type 3 cancer

                    type 1 cancer      type 2 cancer

Another visualisation approach:



Example of using more than 1 feature to decide whether tumor is malignant (M) or benign (B)



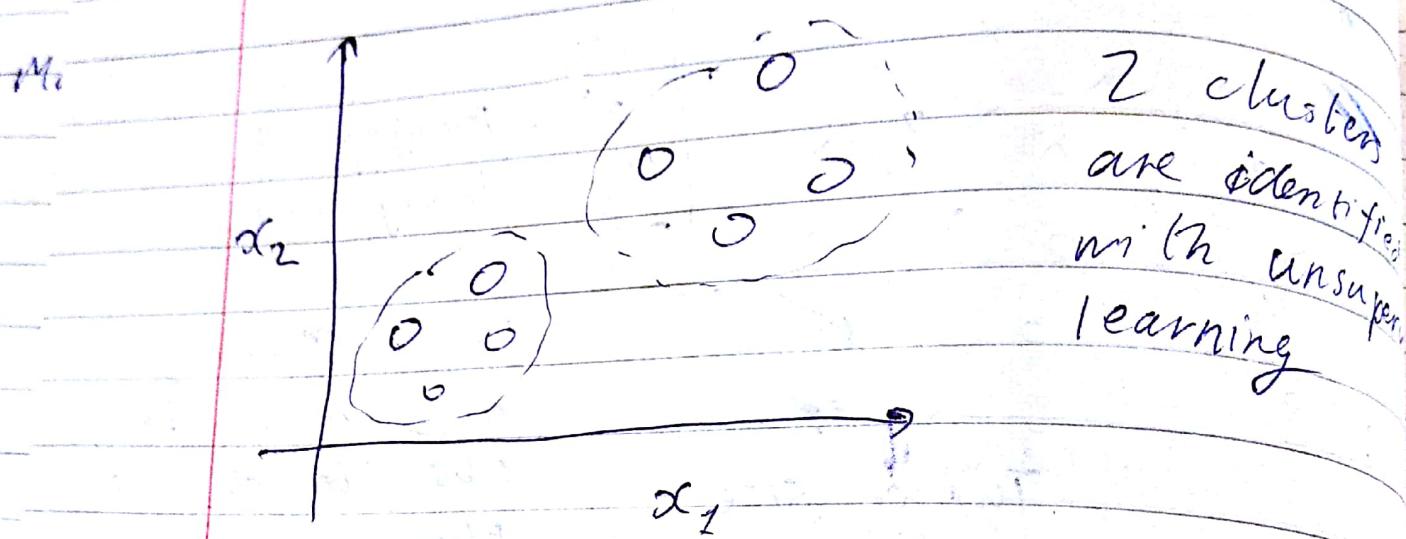
Key :  
x - Malignant (M)  
o - Benign (B)

The learning algo must be able to decide where the boundary lies between M & B

Based on certain values of age & tumor size we can find the probability of the tumor being malignant or benign

Most machine learning problems have many more features. Support vector machines allow you to deal with infinite features with a mathematical trick

Unsupervised learning - data does not contain labels to be predicted (unlike supervised learning). Main objective is to find structure in the dataset



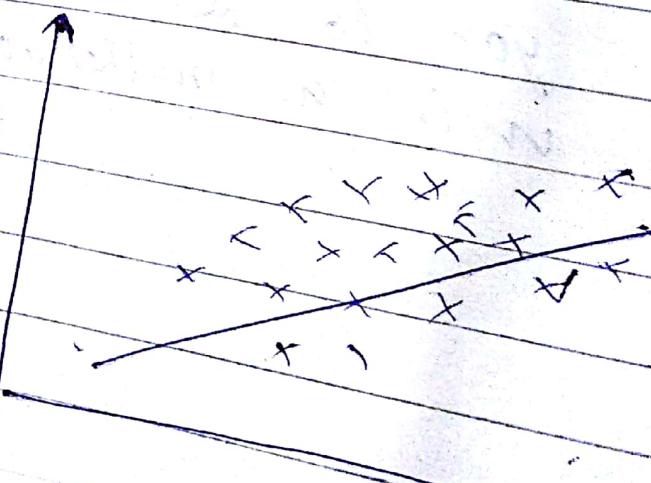
Clustering problem

Linear regression

A model for modeling the relationship between a numerical variable  $y$  & one or more explanatory variables (or features) as a linear function.

~~regression problems usually~~

example  
of supervised  
learning



Important notation:

$m$  - number of training examples

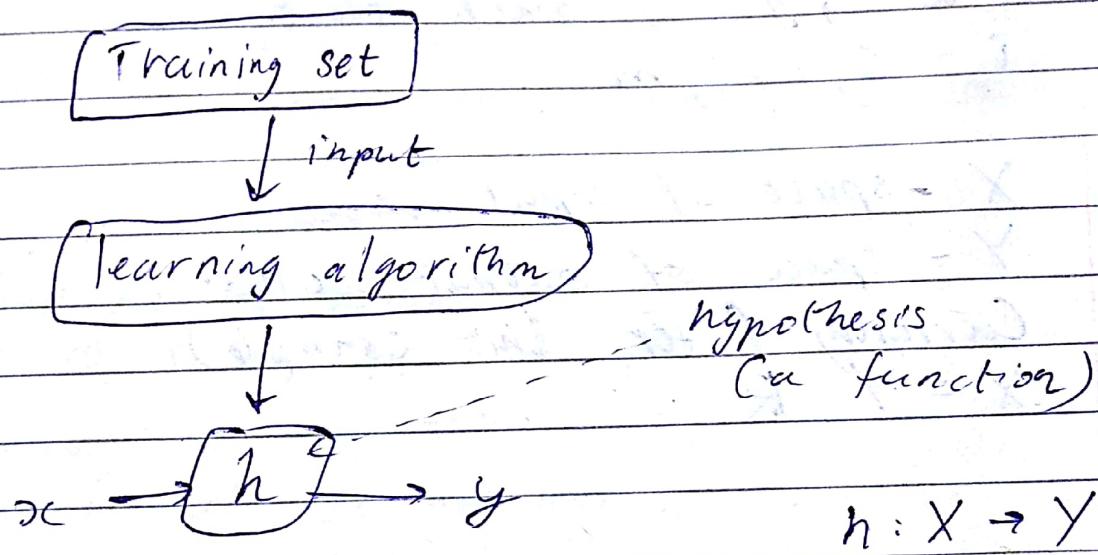
$x$  - input variable / feature

$y$  - output variable / label or target variable

$(x, y)$  - single training example

(a single row in dataset)

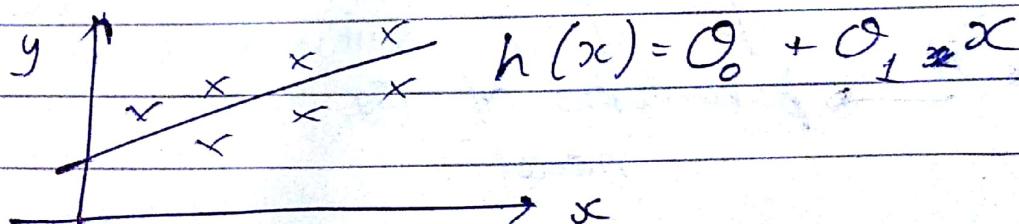
$(x^{(i)}, y^{(i)})$  -  $i$ th training example  
( $i$ th row in dataset)



$h$  maps from  $x$ 's to  $y$ 's

Representing  $h$ : Only for the current time being

$$h_0(x) = \theta_0 + \theta_1 x \quad \} \text{ linear function}$$



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

↳ Linear regression with one variable or univariate linear regression ( $x$  is the variable)

The entire data set:

A list<sup>with</sup>  $m$  training examples

$(x^{(i)}, y^{(i)})$  such that

$$i = 1, 2, \dots, m$$

$X$  - space of input values

$Y$  - space of output values

Currently (for this example)

$$X = Y = \mathbb{R}$$

log  
pro:  
usy

Cost function to help figure out ~~the~~ how fit the best possible line for data

to

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

{ How do we  
choose these  
parameters? }

→ parameters of the model

predicted  
value

Main idea of a cost function:

choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training examples  $(x, y)$

Formally:

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

makes math & calculations easier

means:

find me the values of  $\theta_0$  &  $\theta_1$  that causes this expression to be at its minimum

To add some more notation:

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

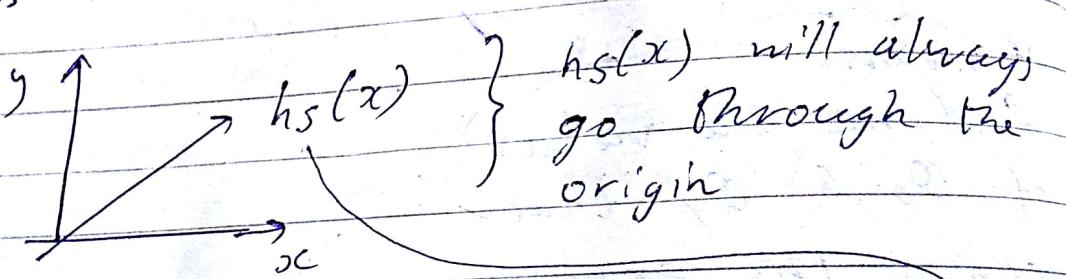
cost  
function/  
squared  
error func

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

- The accuracy of a hypothesis function is measured with a cost function
- The cost function is an average difference thus it is divided by m
- As a convenience for the computation of gradient descent the cost function is multiplied with  $\frac{1}{2}$

Assume simplified  $h_0(x)$  for now:

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



and goal is : minimise  $J(\theta_1)$

s for simplified

$$h_s(x)$$

$$J(\theta_1)$$

for fixed  $\theta_1$ , this

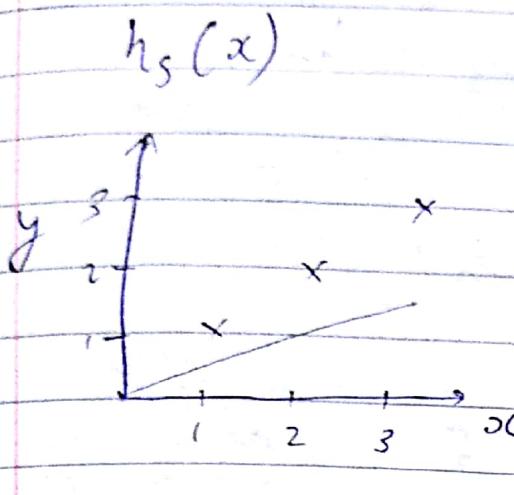
is a function of  $x$

function of

parameter  $\theta_1$

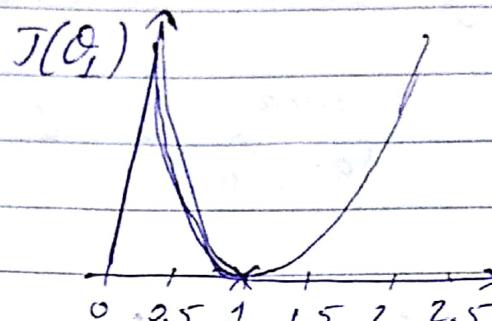
$\theta$  controls the slope of line

The question rests: why is  $J(\theta)$  an average rather than an overall aggregate of cost? Reason is such that if 2 models are to be compared we don't want the # of points to increase the overall cost hence average is used (compare model fit across datasets)



$$\theta_1 = 0.5$$

$$J(\theta_1)$$



the training set & the test set can be divided into different proportion hence

there must be a metric that must measure cost independent of data points

$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] = \frac{1}{6} (3.5) \approx 0.58$$

The minimum of  $J(\theta_1)$  occurs at  $\theta_1 = 1$

### Gradient descent

Will help us minimise the value of  $J(\theta_1, \theta_2)$  by changing values of  $\theta_1, \theta_2$

✓ used in other places as well

### Outline of procedure:

- start with some  $\theta_1, \theta_2$
- keep changing  $\theta_1, \theta_2$  to reduce  $J(\theta_1, \theta_2)$  until we hopefully end up at ~~is~~ minimum

$a := b \rightarrow$  true or false  
 $a := b \rightarrow$  assign value of  $b$  to  $a$

Gradient descent is actually an algo that can be applied on any such problem:

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

repeat until convergence {

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

(for  $j = 0 \text{ } \& \text{ } j = 1$ )

learning rate ( $\alpha$ )  
controls big the are)

}

$\theta_0 \text{ } \& \text{ } \theta_1$  must be updated simultaneously. Serial Implementation:

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Incorrect simultaneous implementation:

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

~~$\theta_1$~~

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

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

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

If updates are not simultaneous then the algorithm won't be called gradient descent & it may behave in strange ~~ways~~ unexpected ways since there is a change in the algorithm

The issue with gradient descent is that there might ~~be~~ be many various local minima & the question of which local minima you will arrive to depends on the initial values of  $\theta_0$  &  $\theta_1$

$$\text{repeat until convergence} \quad \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right.$$

3

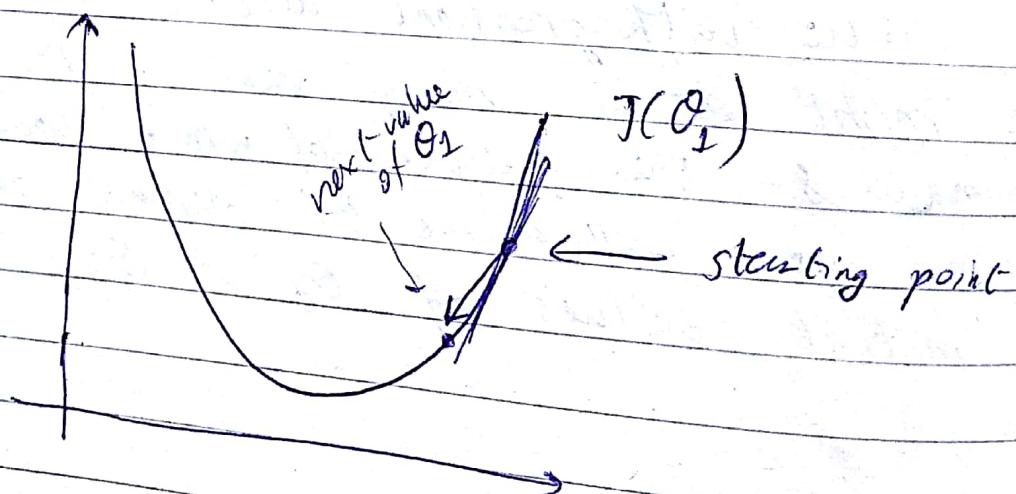
derivative

learning  
rate

(determines step how big  
the update is on each  
iteration)

Suppose simpler cost function:

$$\min_{\theta_1} J(\theta_1) \quad \text{where } \theta_1 \in \mathbb{R}$$



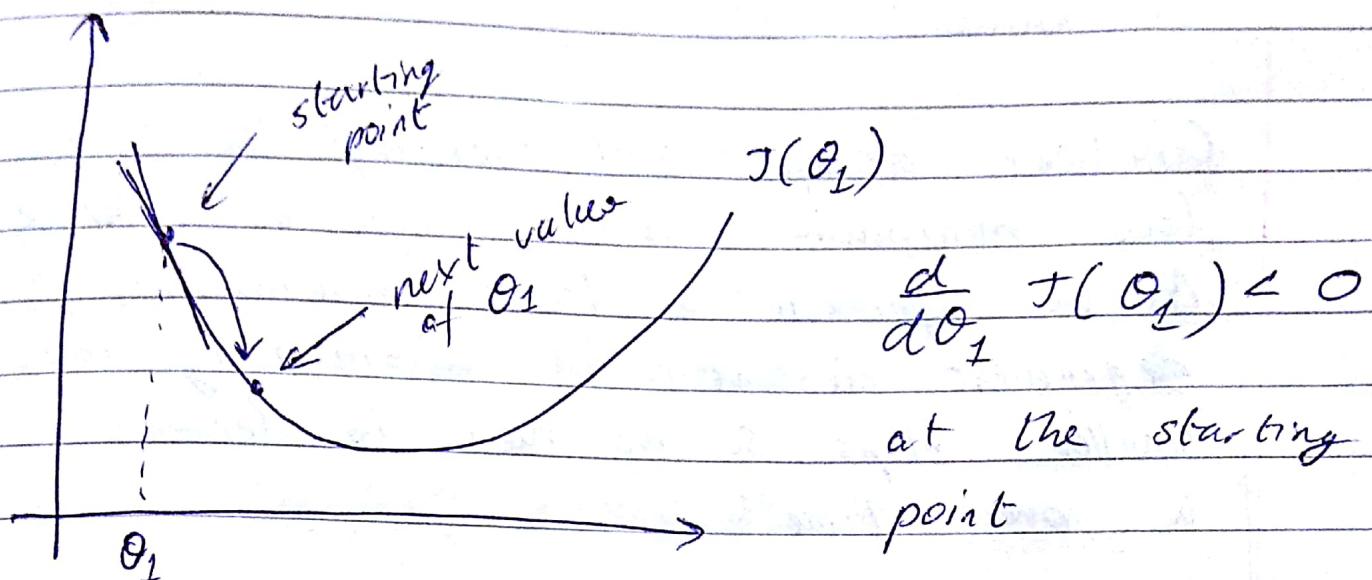
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

calculates

$\frac{d}{d\theta_1} J(\theta_1) > 0$  slope at a  
particular value  
at starting point hence

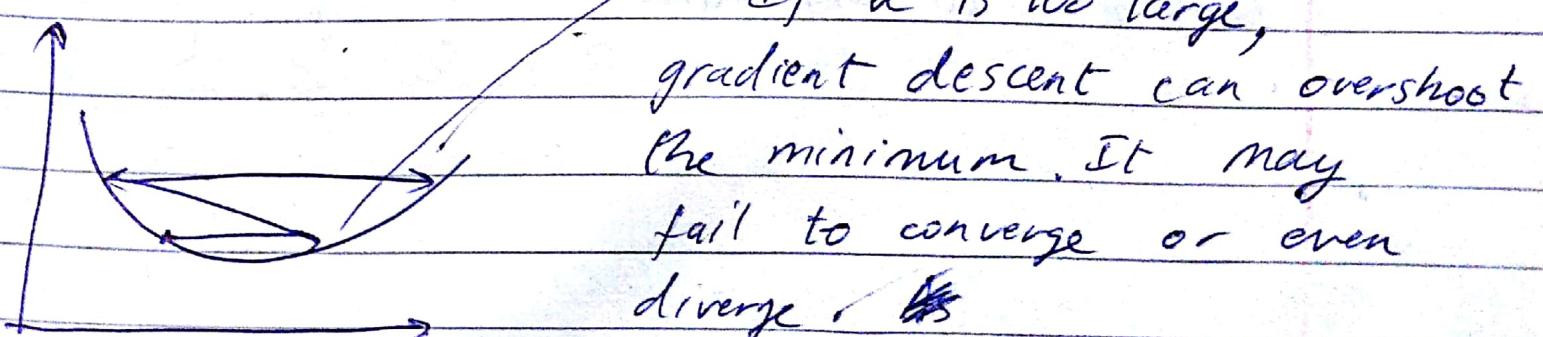
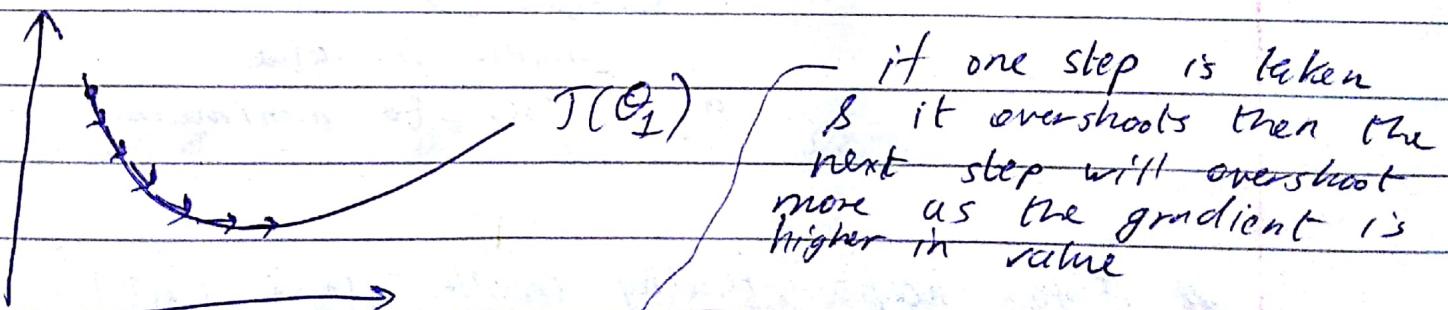
$\theta_1 := \theta_1 - \alpha (+ve)$ , thus  $\theta_1$  decreases in value

Alternatively:



$\theta_1 := \theta_1 - \alpha (-ve)$  thus  $\theta_1$  increases in value

If  $\alpha$  is too small, gradient descent can be slow



If the initial value of  $J(\theta_1)$  reaches a local minimum then the gradient descent algo will stop as at the minimum  $\frac{d}{d\theta_1} J(\theta_1) \geq 0$  hence no more change in value.

Gradient descent will converge to local minimum with fixed value of  $\alpha$ . As we approach a local minimum, gradient descent will automatically take smaller steps. So no need to decrease  $\alpha$  over time.

→ will take smaller steps due to decreasing slope hence decrease or increase will be smaller in value.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$\underbrace{\hspace{10em}}$   
smaller in value  
as closer to minimum

Algo keeps taking smaller steps until converges to minimum

## Combining gradient descent & linear regression model

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

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

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

Linear regression

Gradient descent algo

repeat until convergence {

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

for  $j = 1 \dots n$

}

Finding out value of slope:

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

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

$$\text{For } j=0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) =$$

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

$$\text{For } j=1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) =$$

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

Gradient descent for linear regression

repeat until convergence {

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

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

}

update  $\theta_0$  &  $\theta_1$  simultaneously

Issue with gradient descent:

- susceptible to local optima

simply

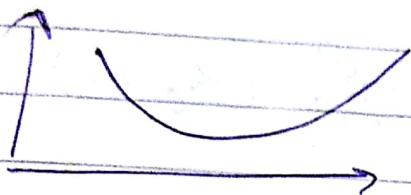
$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

↳ different minimum

value depending on how

$\theta_0$  &  $\theta_1$  are initialised

The cost function for linear regression always is a bow shaped function (a convex function)



} convex - always has 1 single global optimum

Nence gradient descent will always converge to the global optimum for linear regression

Gradient descent algo also called "batch" gradient descent

↳ batch because each step gradient descent uses all the training examples

There are other different versions of existing gradient descent

Finding out the value of the slope  
(General method):

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

$$= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^n \theta_i x_i - y \right)$$

$$= (h_{\theta}(x) - y) (x_j)$$

## Linear Algebra

Matrix : rectangular array of numbers

Dimension of matrix : row x columns

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

2 x 3 matrix

- or

$\mathbb{R}^{2 \times 3}$  matrix

Since A is a matrix, then:

$A_{ij}$  = "i, j entry" in the  $i$ th row &  
 $j$ th column

so:

$$A_{32} = 1, A_{23} = 6$$

Vector - special type of matrix - an  $n \times 1$  matrix  
example:

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 145 \end{bmatrix} \quad \begin{array}{l} \text{4-dimensional} \\ \text{vector} \end{array} \quad \begin{array}{l} \text{only 1} \\ \text{column} \end{array}$$

or

$$\mathbb{R}^4 \text{ vector}$$

$$y_i = i\text{th element} \quad \text{so } y_1 = 460, y_2 = 232$$

Upper case letters are used to refer to matrices  
Lower case letters are used to refer to scalars & vectors

### Matrix Addition

you can only  
add matrices of  
the same dimension

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0.5 \\ 4 & 10 \\ 3 & 2 \end{bmatrix}$$

$3 \times 2$

$3 \times 2$

$3 \times 2$

## Scalar Multiplication

$$3 \times \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix} =$$
$$3 \times 2 \quad \quad \quad 3 \times 2$$
$$= \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} \times 3$$

$$\begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} / 4 = \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1.5 & 0.75 \end{bmatrix}$$
$$2 \times 2 \quad \quad \quad 2 \times 2 \quad \quad \quad 2 \times 2$$

The dimensions will remain the same

## Vector - Matrix Multiplication

$$\begin{bmatrix} 1 & 3 \\ 4 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 5 \end{bmatrix} = \begin{bmatrix} 16 \\ 4 \\ 7 \end{bmatrix}$$
$$3 \times 2 \quad \quad \quad 2 \times 1 \quad \quad \quad 3 \times 1$$

$$1 \times 1 + 3 \times 5 = 16$$

$$4 \times 1 + 0 \times 5 = 4$$

$$2 \times 1 + 1 \times 5 = 7$$

Suppose the following house size data:

2104  
1416  
1534  
852

$$\text{Hypothesis: } h_0(x) = -40 + 0.25x$$

We need to evaluate ~~the~~ a prediction on each of sizes. This can be done through matrices & vectors:

$$\begin{matrix} 4 \times 2 \\ \begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 \\ 0.25 \end{bmatrix} = \end{matrix}$$

$$= \begin{bmatrix} -40 \times 1 + 0.25 \times 2104 \\ -40 \times 1 + 0.25 \times 1416 \\ -40 \times 1 + 0.25 \times 1534 \\ -40 \times 1 + 0.25 \times 852 \end{bmatrix} \begin{matrix} h_0(2104) \\ h_0(1416) \\ \vdots \end{matrix}$$

4 × 1 matrix

1st way

In code the above operation can be written in a single line as such:

prediction = Data matrix \* parameters

2nd way) Compare this to using a loop:

for  $i = 1: 4$   
predraction(i)

By using the 2st way the code is simplified by a great degree & it is computationally efficient as well

### Matrix - Matrix multiplication

$$\begin{matrix} 2 \times 3 \\ \left[ \begin{matrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{matrix} \right] \end{matrix} \times \begin{matrix} 3 \times 2 \\ \left[ \begin{matrix} 1 & 3 \\ 0 & 1 \\ 5 & 2 \end{matrix} \right] \end{matrix} = \begin{matrix} 2 \times 2 \\ \left[ \begin{matrix} 11 & 10 \\ 9 & 14 \end{matrix} \right] \end{matrix}$$

$$\left[ \begin{matrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{matrix} \right] \times \left[ \begin{matrix} 1 \\ 0 \\ 5 \end{matrix} \right] = \left[ \begin{matrix} 11 \\ 9 \end{matrix} \right]$$

$$\left[ \begin{matrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{matrix} \right] \times \left[ \begin{matrix} 3 \\ 1 \\ 2 \end{matrix} \right] = \left[ \begin{matrix} 10 \\ 14 \end{matrix} \right]$$

Hence:

$$\left[ \begin{matrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{matrix} \right] \times \left[ \begin{matrix} 1 & 3 \\ 0 & 1 \\ 5 & 2 \end{matrix} \right] = \left[ \begin{matrix} 11 & 10 \\ 9 & 14 \end{matrix} \right]$$

## Application:

Suppose we have the following houses & we wish to compute ~~the~~ its predictions:

2104, 1416, 1534, 852

3 competing hypothesis are available:

$$1) h_0(x) = -40 + 0.25x$$

$$2) h_0(x) = 200 + 0.1x$$

$$3) h_0(x) = -150 + 0.4x$$

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 & 200 & -150 \\ 0.25 & 0.1 & 0.4 \end{bmatrix} =$$

An efficient  
manner of  
applying 3  
hypothesis on 4  
data points &  
can be written  
with 1 line of  
code as well

$$= \begin{bmatrix} 486 & 910 & 692 \\ 314 & 342 & 416 \\ 344 & 353 & 464 \\ 173 & 285 & 191 \end{bmatrix}$$

each column is  
a list of predictions  
from one hypothesis

2nd hypothesis  
predictions

First hypothesis  
predictions

3rd hypothesis  
predictions

## Matrix Multiplication properties

$$A \times B \neq B \times A$$

if  $A$  &  $B$  are matrices i.e.  
matrix multiplication is not  
commutative

Matrix multiplication is associative:

## Identity matrix:

Denoted by  $I$  or  $I_{n \times n}$

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

For any matrix A:

$$A \times I = I \times A = A$$

$m \times n \quad n \times n \quad m \times m \quad m \times n \quad m \times n$

↑      ↑

different

because if  $A$  is  $n \times m$  then:  
 $A A^{-1} = I_n$  &  $A^{-1}A = I_m$   
 $\Rightarrow I_n \neq I_m$

## Matrix inverse

must be  
a square matrix

If  $A$  is an  $m \times m$  matrix and if it has an inverse,

$$AA^{-1} = A^{-1}A = I$$

$(n \times m) (m \times n)$        $(m \times n) (n \times m)$

E.g. 
$$\begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix} \begin{bmatrix} 0.4 & -0.1 \\ -0.05 & 0.075 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_{2 \times 2}$$

A                   $A^{-1}$                    $AA^{-1}$

Some matrices do not have an inverse —  
singular matrix or degenerate matrix

e.g.: 
$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

Matrix Transpose convert rows  
into columns

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix}_{2 \times 3}, \quad A^T = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix}_{3 \times 2}$$

Let  $A$  be a  $m \times n$  matrix, and let  $B = A^T$   
Then  $B$  is an  $n \times m$  matrix and

$$B_{ij} = A_{ji}$$