

Q) What is machine learning?

- field of study that gives computers the ability to learn without being explicitly programmed.
- A computer program is said to 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:

- Supervised learning (features & labels)
- Unsupervised learning (only features)

Others: Reinforcement learning, recommender systems.

* Linear regression with one variable
→ model representation.

$y = mx + c$

Notation:

m = Number of training Example

x 's = "input" Variable/feature

y 's = "Output" Variable / "target" Variable

(x, y) - One training Example

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

Regression

Training Set \rightarrow Continuous Value

(Data to extract) \rightarrow nominal, basic Classification

(at first step) \rightarrow discrete, continuous, etc.

Learning Algorithm \rightarrow Discrete Value

divide

test set
(Input Values)



h

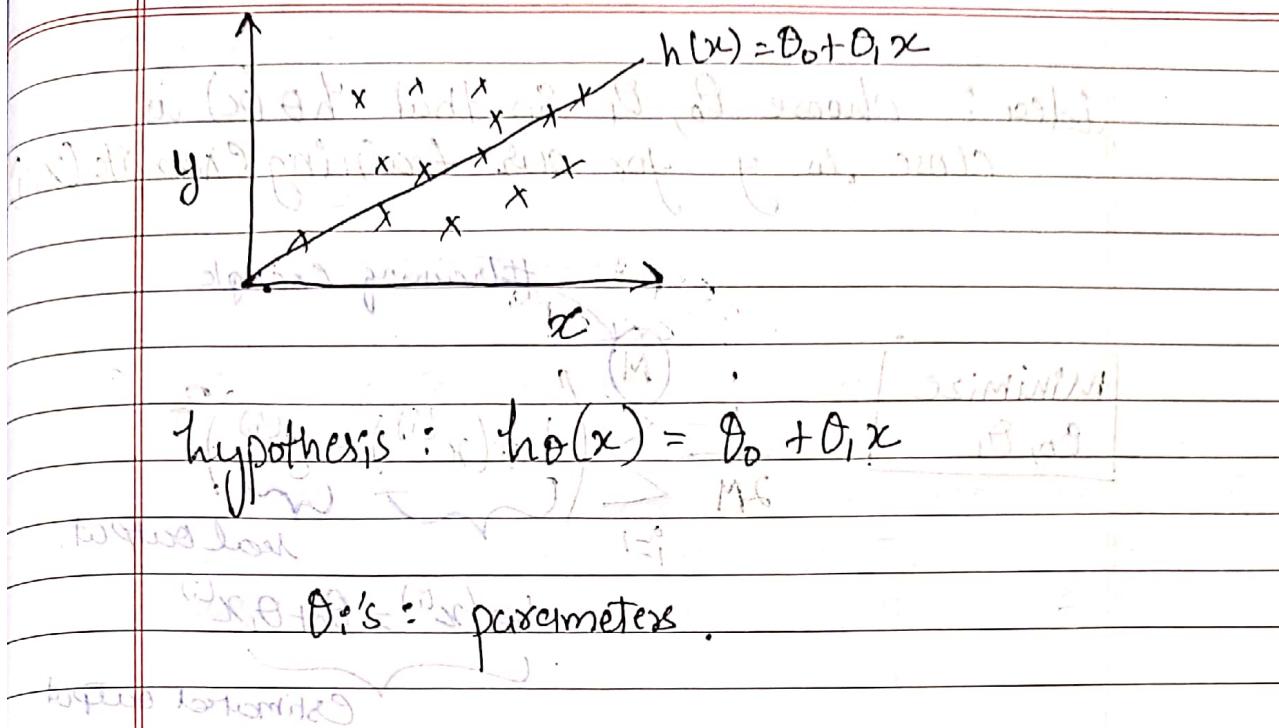


estimated Value
 \hat{y} (output)

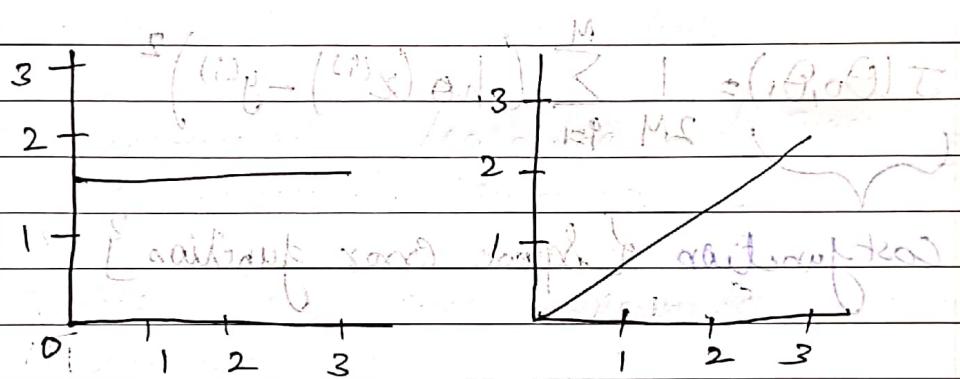
* How do we represent h ?

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

Shorthand : $h(x)$



#



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

By fitting line

$$\theta_0 = 1 \quad \theta_1 = 0.5$$

$$\left(\frac{\partial J}{\partial \theta_0} - \frac{\partial J}{\partial \theta_1} \right)_{\theta_0=1, \theta_1=0.5} \leq 0$$

Idea: choose θ_0, θ_1 so that $h_\theta(x)$ is close to y for our training example $(x^{(i)}, y^{(i)})$

Training Example

Minimize
 θ_0, θ_1

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

real output
 $h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$

Estimated output

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

Cost function of square error function

Cost function

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

parameters: θ_0, θ_1

Cost function:

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

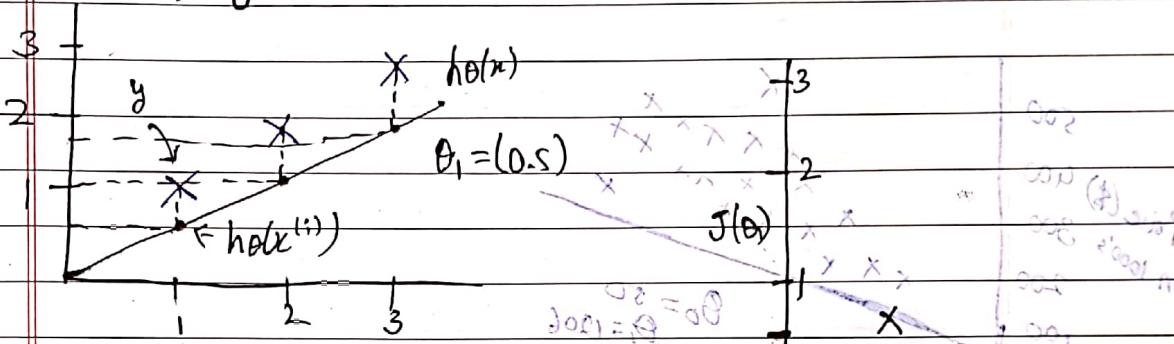
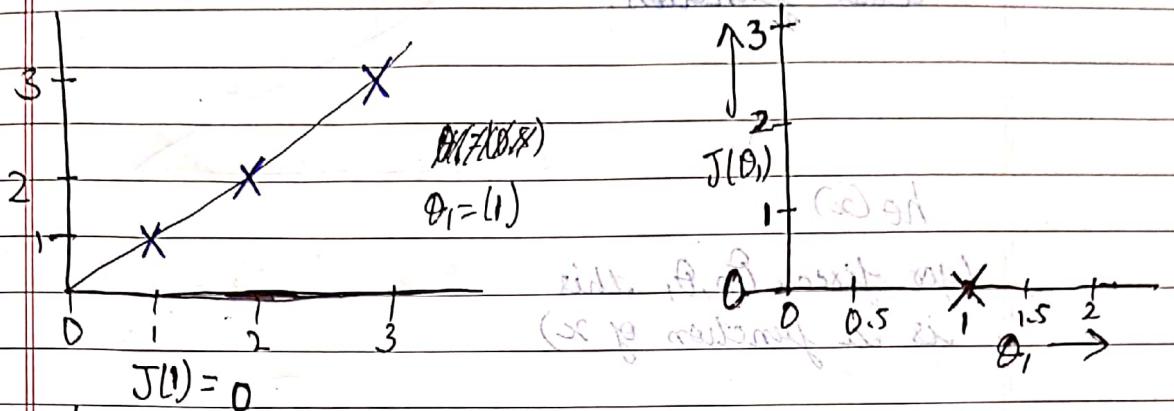
Goal: minimum $J(\theta_0, \theta_1)$

Simplified

$$h_\theta(x) = \theta_1 x ; \theta_0 = 0$$

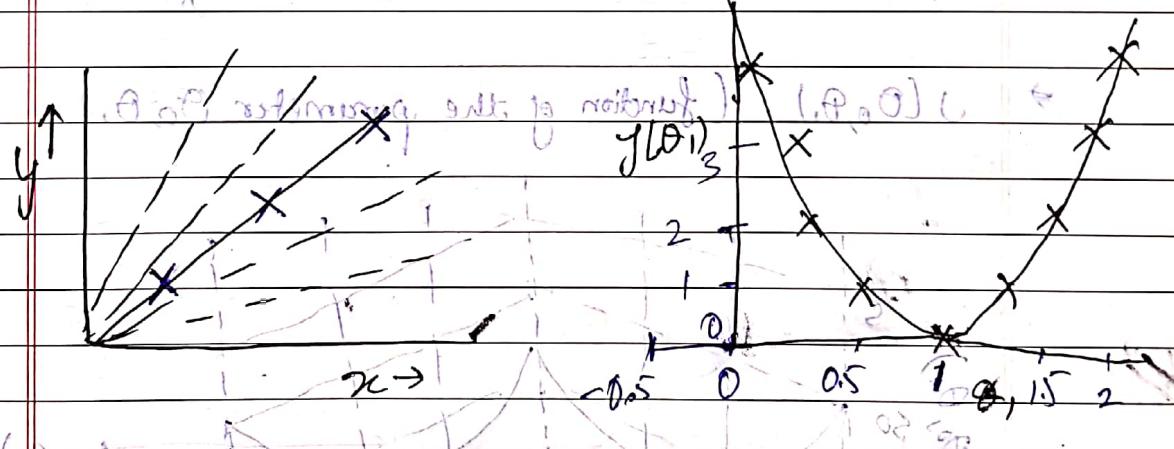
$$J(\theta_1) = \frac{1}{2M} \sum_{i=1}^M (h_\theta(x^{(i)}) - y^{(i)})^2$$

get min $\theta_0(\theta_1)$, normalise $J(\theta_1)$ w.r.t θ_1



$$J(0.5) = \frac{1}{2M} [(0.5-1)^2 + (r_2)^2 + (1.5-3)^2]$$

$$= \frac{1}{2M} (3.5) = 0.58$$



(θ_0, θ_1)

sof

15

10

5

0

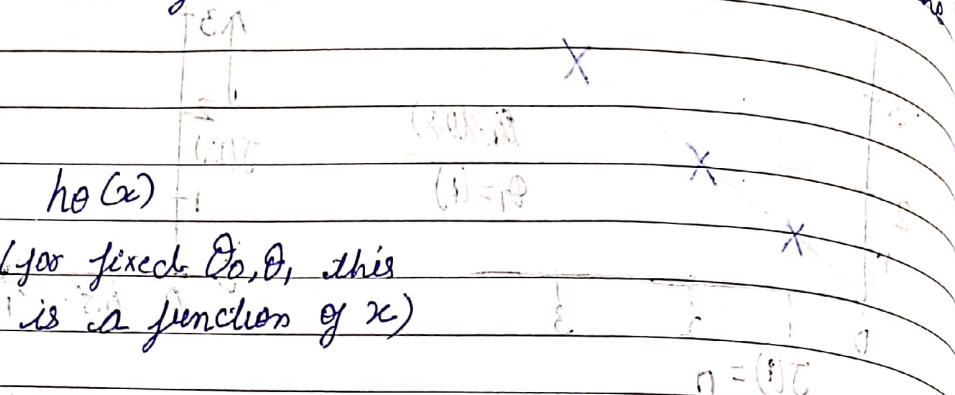
-5

-10

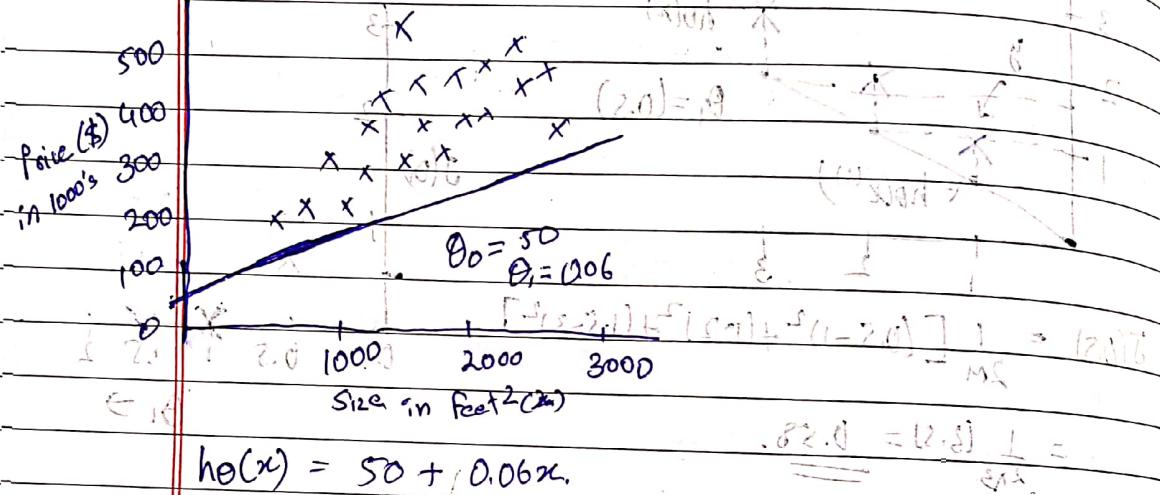
-15

linear

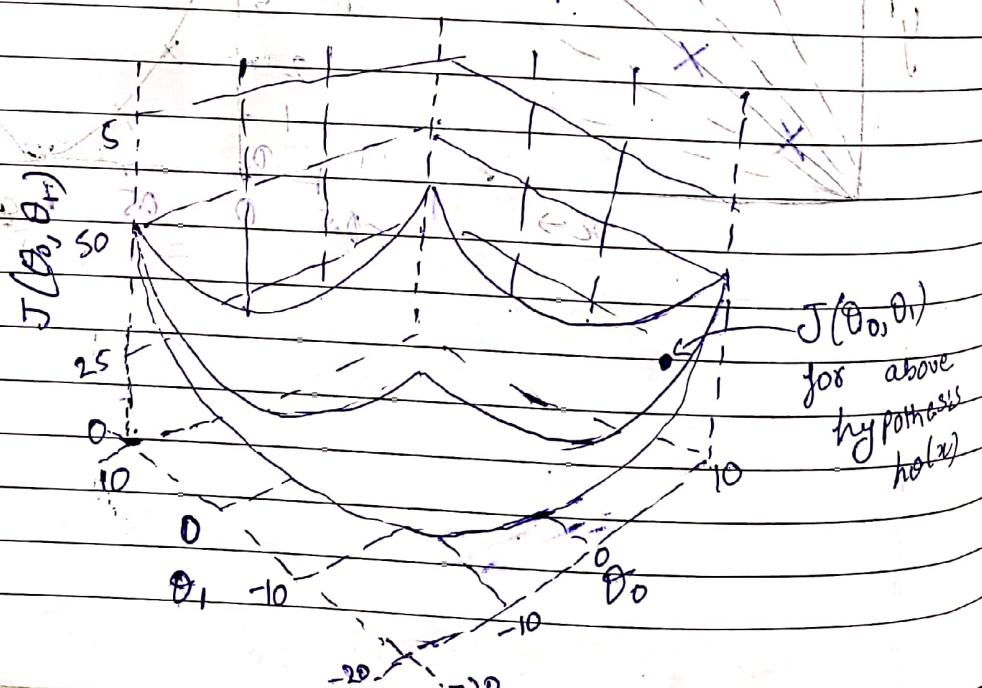
⇒ Better the regression line minimum the cost function.



(for fixed θ_0, θ_1 this is a function of x)



⇒ $J(\theta_0, \theta_1)$ (function of the parameter θ_0, θ_1)

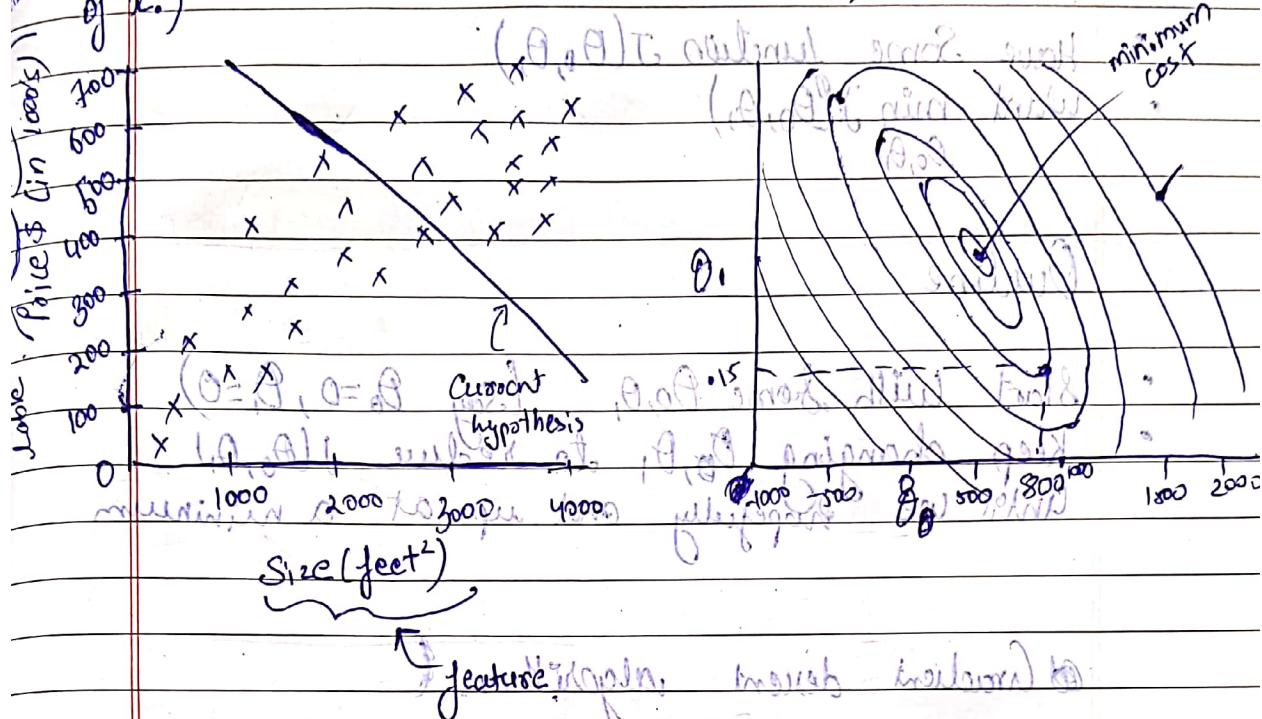


$h(x)$

(for fixed θ_0, θ_1 , this is a function of x)

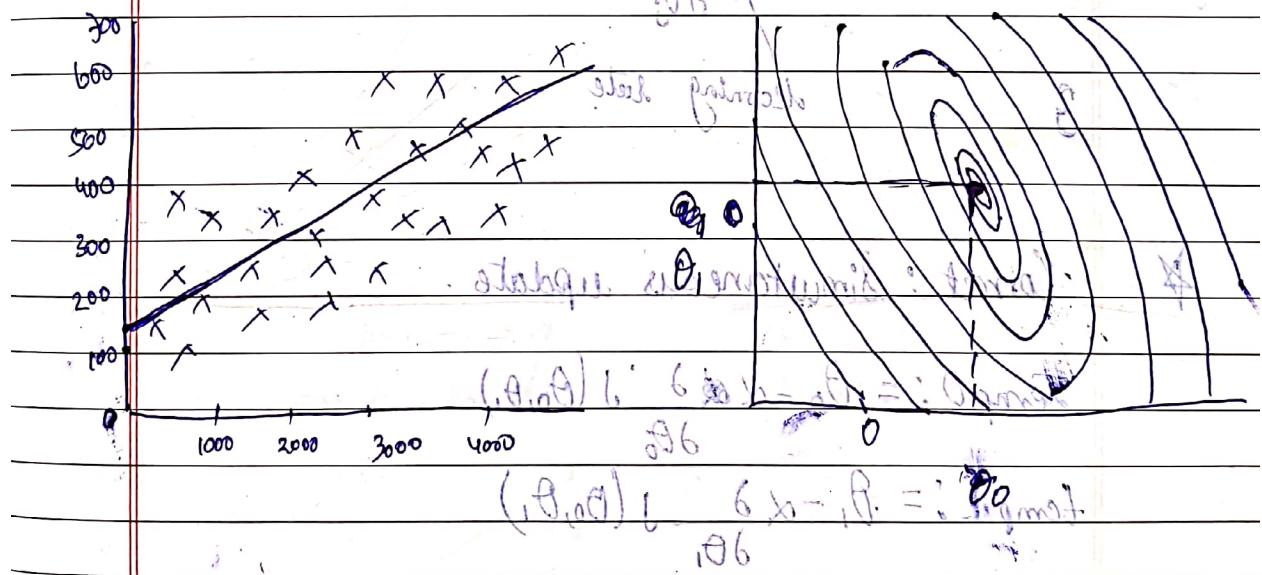
 $J(\theta_0, \theta_1)$

(function of the parameter θ_0, θ_1)

 $h_\theta(x)$

(for fixed θ_0, θ_1 , this is a

(function of x) $\hat{y} = \theta_0 + \theta_1 x$



$$\theta_0 \text{ const} = 0.00$$

$$\theta_1 \text{ const} = 0.00$$

Gradient descent

Have some function $J(\theta_0, \theta_1)$

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

Outline

- Start with some θ_0, θ_1 (say $\theta_0 = 0, \theta_1 = 0$)
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
 Until we hopefully end up at a minimum
(steepest descent)

② Gradient descent algorithm

repeat until convergence

2

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

3

decreasing rate

* Correct: simultaneous update.

$$\text{temp}0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp}1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}0$$

$$\theta_1 := \text{temp}1$$

* Incorrect :

$$(1) \theta_0 \text{ (initial)} = \theta_0$$

196

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

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

local opt is N.F.

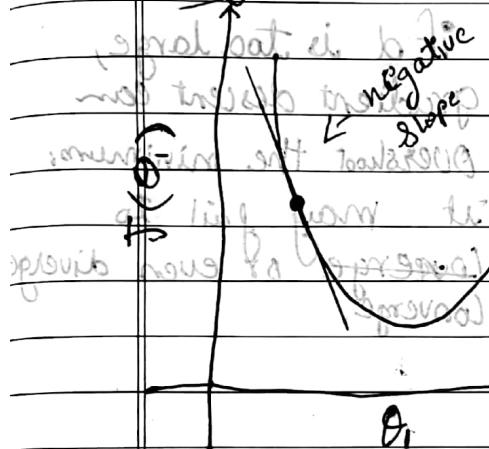
~~$\text{temp}_1 := \theta_0 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$~~

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

*

Why minus (-) in

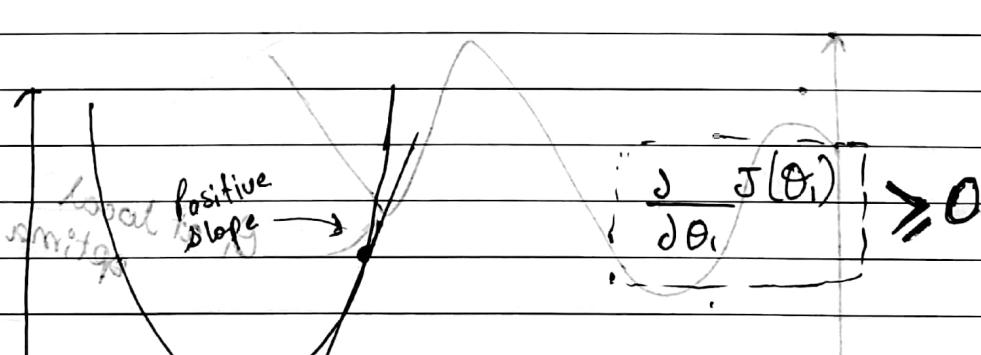
$$\theta_i := \theta_i - \alpha \frac{\partial}{\partial \theta_i} J(\theta)$$



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

$$\left| \frac{\partial}{\partial \theta_i} J(\theta_i) \right| \leq 0$$

$$\theta_i := \theta_i - \alpha (\text{negative number})$$



$$\theta_i := \theta_i - \alpha (\text{Positive number})$$

(197)

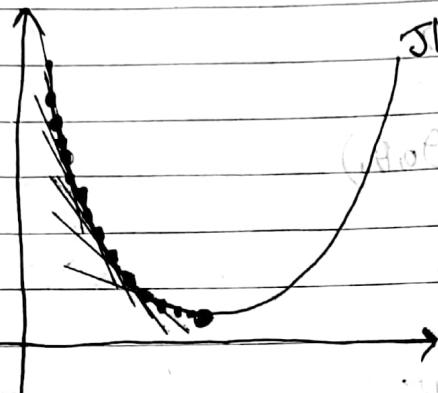
$$\theta_0 = \theta$$

196

increasing
value

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

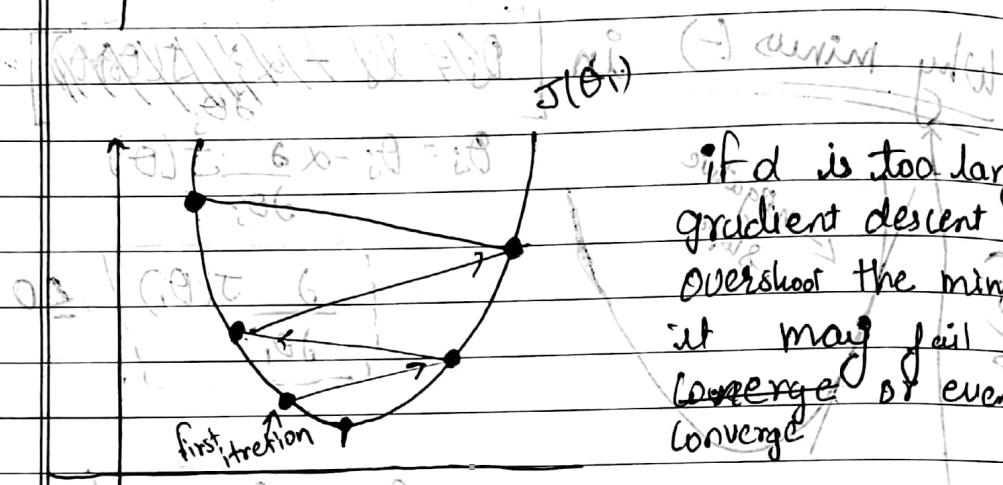
(θ_1, θ_2) & $\theta_1 = \theta_2 = 0$ =: origin



$J(\theta)$

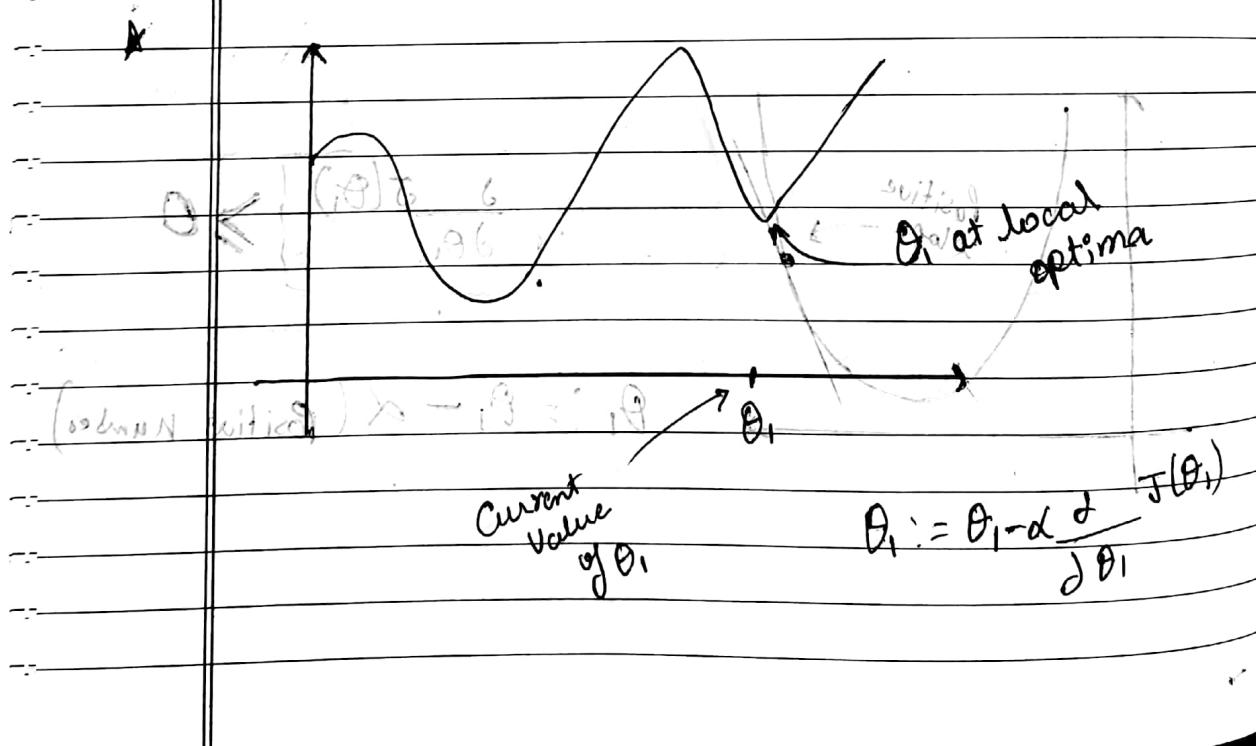
if α is too small,
gradient descent can
be slow.

1 step = 10



$J(\theta)$

if α is too large,
gradient descent can
overshoot the minimum,
it may fail to
converge or even diverge

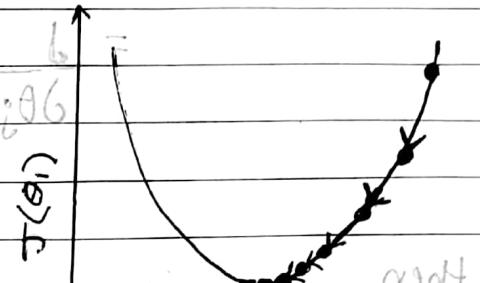


Current
value
of θ_1

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

* gradient descent can converge to a local minimum, even with the learning rate α fixed.

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



As we approach a local minimum, gradient

descent will automatically take smaller steps. So, no

need to decrease α over time.

bigger the slope $\left(\frac{\partial J(\theta)}{\partial \theta_i} \right)$ larger the step.

* linear regression model (one variable)

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

(cost function) $J(\theta)$

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

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

$$(0 - h_\theta(x^{(i)}) + \theta_j x_j^{(i)}), \left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right)$$

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

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

then,

$$j=0 : \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot \frac{\partial}{\partial \theta_0} (\theta_0 + \theta_1 x^{(i)} - y^{(i)})$$

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

$$= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot (1 + 0 - 0) = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})$$

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

$j=1$ (feature)

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

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

$$= \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot (0 + x^{(i)} - 0)$$

$$\begin{aligned}
 &= \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot x^{(i)} \\
 &= \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}
 \end{aligned}$$

rumus agaknya, cari θ iaitu nilai "parameter" = θ

* If repeat until convergence

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

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

update θ_0 and θ_1 simultaneously

* "Batch"

→ Each step of gradient descent uses all the training examples.

MATRICES AND VECTORS.

Matrices: Rectangular arrays of numbers.

Dimension of matrix: number of rows × number of columns.

Example:	$\begin{bmatrix} 1 & 0 & 2 \\ 3 & 2 & 1 \\ 4 & 0 & 2 \end{bmatrix}$	$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$	A	3×3	9 elements
.	$\begin{bmatrix} 1 & 0 & 2 \\ 3 & 2 & 1 \\ 4 & 0 & 2 \end{bmatrix}$	$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$	4×2 matrix	4×2	8 elements

$$A = \begin{bmatrix} 1402 & 8102 \\ 2047 & 1112 \\ 9234 & 378 \\ 071 & 9279 \end{bmatrix}$$

A_{ij} = "i,j entry" in the i^{th} row, j^{th} column.

$$A_{22} = 1112, A_{12} = 8102, A_{41} = 071$$

* Vector : An $n \times 1$ matrix

$$\text{Ex. } (10 \times 1) \times 1 = 10$$

$$y = \begin{bmatrix} 460 \\ 320 \\ 123 \\ 145 \end{bmatrix}$$

4-Dimensional Vector.

y_i = i^{th} element

Up to now we have learned about 2D and 3D vectors.

$$y_1 = 460, y_2 = 320, y_3 = 123, y_4 = 145 \text{ etc.}$$

⇒ 1-indexed vs 0-indexed

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

$$y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

M

Y

Y

Y

for example a 2D coordinate system is defined as follows:

⇒ We use A, B, C, X for matrices or

systems of equations

and, a, b, x, y for vectors.

* matrix addition

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

xistum 1×8

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix}_{3 \times 2} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \end{bmatrix}_{2 \times 2} = \text{error}$$

dimensions are not same

* Scalar Multiplication.

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

xistum 1×3

scalar (dimensions match) (dimensions not same)

* Scalar division

$$\begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix}_{2 \times 2} \div 4 = \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix}_{2 \times 2} = \begin{bmatrix} 1 & 0 \\ 3/2 & 3/4 \end{bmatrix}_{2 \times 2}$$

qu. ment. into tens (or scalar) for

* Combination of Operands

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}_{1 \times 3} + \begin{bmatrix} 3 & 2 & 1 \\ 2 & 3 & 0 \\ 1 & 0 & 1 \end{bmatrix}_{3 \times 3} = \begin{bmatrix} 3 & 2 & 1 \\ 2 & 3 & 0 \\ 1 & 0 & 1 \end{bmatrix}_{3 \times 3}$$

$$\begin{bmatrix} 3 & 0 & 1 \\ 12 & 0 & 1 \\ 6 & 5 & 1 \end{bmatrix}_{3 \times 3} = \begin{bmatrix} 2 \\ 12 \\ 10 \end{bmatrix}_{3 \times 1} \text{ 3-dimensional vector}$$

$1 \times 3 + 3 \times 3 + 1 \times 1$

* Matrix-Vector multiplication x system

$$\begin{bmatrix} 1 & 3 \\ 4 & 0 \\ 2 & 1 \end{bmatrix}_{3 \times 2} \times \begin{bmatrix} 2 \\ 5 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} (1 \times 2) + (3 \times 5) \\ (4 \times 2) + (0 \times 5) \\ (2 \times 2) + (1 \times 5) \end{bmatrix} = \begin{bmatrix} 16 \\ 8 \\ 7 \end{bmatrix}_{3 \times 1} \text{ matrix}$$

Same, means possible.
dimension of resultant matrix

$$A_{m \times n} \times x_{n \times 1} = y_{m \times 1}$$

A is $m \times n$ matrix (m rows, n columns)
 x is $n \times 1$ vector (n-dimensional vector)

$\begin{bmatrix} 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 \end{bmatrix} + \dots + \begin{bmatrix} 0 & 1 \end{bmatrix}$

To get y_i , multiply A 's i^{th} row with elements of vector x , and add them up.

Example : $\begin{bmatrix} 1 & 2 & 1 & 5 \end{bmatrix}_{3 \times 4} \times \begin{bmatrix} 1 \\ 3 \\ 2 \\ -1 \end{bmatrix}_{4 \times 1} = \begin{bmatrix} 14 \\ 13 \\ -7 \end{bmatrix}_{3 \times 1}$

$$\begin{bmatrix} 1 & 2 & 1 & 5 \end{bmatrix}_{3 \times 4} \times \begin{bmatrix} 1 \\ 3 \\ 2 \\ -1 \end{bmatrix}_{4 \times 1} = \begin{bmatrix} 14 \\ 13 \\ -7 \end{bmatrix}_{3 \times 1}$$

x system

$$\begin{bmatrix} 1x1 + 2x3 + 1x2 + 5x(-1) \\ 0x1 + 3x3 + 0x2 + 4x(-1) \\ -1x1 + (-2)x3 + 0x2 + 0x(-1) \end{bmatrix} = \begin{bmatrix} 14 \\ 13 \\ -7 \end{bmatrix}$$

House Sizes (x) $\theta_0 : 2104, 1534, 852$

$$h_{\theta}(x) = -40 + 0.25x$$

$$\begin{array}{c} 2104 = \text{B} \\ 1416 \\ 1534 \\ 852 \end{array} \quad \begin{array}{c} x \\ \theta = \begin{bmatrix} -40 \\ 0.25 \end{bmatrix} \\ \text{A} \end{array}$$

$\text{for } x = \begin{bmatrix} 2104 \\ 1416 \\ 1534 \\ 852 \end{bmatrix}$
 $\text{for } \theta = \begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$
 $\text{for } h_{\theta}(x) = \begin{bmatrix} (-40 \times 2104) + (0.25 \times 2104) \\ \dots \\ (-40 \times 852) + (0.25 \times 852) \end{bmatrix}$

Brackets in x is 4×1 for matrix 2×1 for θ

θ is 2×1 for θ is 1×2 for matrix 2×1 for $h_{\theta}(x)$

Prediction = Data matrix * parameters (θ)
 (x_0, x_1, \dots, x_n)

★ Matrix - Matrix multiplication

(x) result result

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 3 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 11 & 10 \\ 9 & 14 \end{bmatrix}$$

$2 \times 3 \times 2 = 2 \times 2$ 2×2

$2 \times 1 \times 2 = 2 \times 1$ 2×1

$2 \times 1 \times 1 = 2 \times 1$ 2×1

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 11 \\ 9 \end{bmatrix}$$

x result

x result

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 11 \\ 9 \end{bmatrix}$$

11 9

Matrix multiplication

11 9

Details:

(c) basic form

$$A \times B = P \times C$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

m x n matrix
(m rows, n columns)

n x p matrix
(n rows, p columns)

m x p matrix
(m rows, p columns)

The i^{th} column of the matrix C is obtained by multiplying A with the i^{th} column of B .

~~For example: $x_1 \text{st row of } A \times x_1 \text{st column of } B = \text{first row of } C$~~

(for $i=1, 2, \dots, o$)

A

House sizes (x)

$$\begin{bmatrix} 2104 \\ 1416 \\ 1534 \\ 852 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{Have 3 competing hypotheses:}$$

$$1. h_{\theta}(x) = -40 + 0.25x$$

$$2. h_{\theta}(x) = 200 + 0.12x$$

$$3. h_{\theta}(x) = -150 + 0.4x$$

Matrix
for θ_0 for θ_1

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$$

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

$$\begin{bmatrix} 200 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\begin{bmatrix} 486 \\ 314.8 \\ 344.0 \\ 173 \end{bmatrix} \begin{bmatrix} 1410 \\ 342 \\ 353 \\ 285 \end{bmatrix} \begin{bmatrix} 692 \\ 416 \\ 464 \\ 191 \end{bmatrix}$$

Prediction
of 1st ho Prediction
of 2nd ho

Matrix multiplication properties

* Let A and B be matrices. Then in general
~~AXB ≠ BXA~~ (note commutative) ~~in fact, B~~

Ex. If $A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$, $B = \begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix}$ then $A \times B = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$ i.e. $A \times B \leftarrow mxn \times n \times m$
 $A \times B$ is $m \times m$

Or, $\begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 2 & 2 \end{bmatrix} = B \times A$ is $n \times n$

which is I said (return answer) $AT = A$ ~~exist~~

* Let A and B and C be matrices.

~~if $(A \times B) \times C = A \times (B \times C)$ (Associative)~~
~~"strongly" to "regular"~~

Identity Matrix

Denoted I (or $I_{n \times n}$)

scattered exist

E.g. Example: $A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = A$
1×3	3×3	3×3

then known exist 1×3 $\rightarrow [1 \ 0 \ 0]$ ~~1x3~~

know, x $n \times n$ and in \mathbb{R} matrix $nxn = I_n$

for any matrix A,

$A \cdot I = I \cdot A = A$

$m \times n$	$n \times n$	$m \times m$	$m \times n$
--------------	--------------	--------------	--------------

note: I_n

$A \times B \neq B \times A$
 $AI = IA$

$P = \text{diag}$

$P = \text{diag}$



Inverse and transpose

matrix inverse. contains sd I b/w A & A⁻¹

→ Not all numbers have an inverse.

→ If A is an m x m matrix, and if it has an inverse,

m x m i.e. 8 x 8

$$\rightarrow A(A^{-1}) = \cancel{A} (A^{-1}) A = I$$

→ A = TR^{m x m} (square matrix) because I is square matrix.

→ (matrices) that don't have an inverse are "singular" or "degenerate".

Matrix Transpose

Example:

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \\ 3 & 0 & 9 \end{bmatrix} \quad A^T = A' = \begin{bmatrix} 1 & 0 & 3 \\ 2 & 0 & 0 \\ 0 & 1 & 9 \end{bmatrix}$$

let A be an m x n matrix, and let B = A^T. Then B is an n x m matrix, and

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

$$\Rightarrow A \times B_{12} = A_2 A = 2$$

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

$$B_{32} = 9$$

$$A = A \cdot I = I \cdot A$$

$$A_{22} \cdot A_{23} = 9$$

★ LINEAR REGRESSION WITH

MULTIPLE VARIABLES (Features)

Size(feet ²) (x_1)	No. of bedrooms (x_2)	No. of floors (x_3)	Age of home (years) (x_4)	price (\$1000) (y)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
...

... $\rightarrow \hat{y} = \theta^T x$

Notation :

m = number of training Examples.

n = number of features.

$x^{(i)}$ = input (values) to i^{th} Training Example.

$x_j^{(i)}$ = value of feature j in i^{th} Training example

Hypothesis: $\hat{y} = \theta^T x$

$$\hat{y}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

for convenience of notation, define $x_0^{(i)} = 1$

$$(\theta)^T \cdot x - \theta = 0$$

($\theta_0 = 0$ for first feature (bias))

for convenience of notation, define $x_0 = 1$

$$\text{then } X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

using
 $\theta^T = [\theta_0, \theta_1, \theta_2, \dots, \theta_n]$
 $x^T = [x_0, x_1, x_2, \dots, x_n]$
 $\theta^T x^T = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

then, $\theta^T = [\theta_0, \theta_1, \theta_2, \dots, \theta_n] * x^T = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$h_{\theta}(x) = \theta^T x$ ← Hypothesis.

: wait for it

number of points for training = m

number of features = n

Gradient descent (for multiple Variables)

Cost function :

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

$$x_0\theta_0 + \dots + x_n\theta_n + \theta_0 + \theta_1 + \dots + \theta_n$$

Gradient descent : for minimum cost

Repeat {

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

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

Repeat {

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

(Simultaneously update θ_j for $j=0, \dots, n$)

3

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

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

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

for better understanding:

$$(h_\theta(x^{(i)}) - y^{(i)}) = \left[\begin{array}{c} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \end{array} \right] - \left[\begin{array}{c} 1 \\ x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \end{array} \right] = \left[\begin{array}{c} 0 \\ x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \end{array} \right]$$

$$\theta = \left[\begin{array}{c} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{array} \right] \in \mathbb{R}^{n+1}$$

$$(n+1) \times 1 \quad \rightarrow \quad m \times (n+1)$$

$$(h_{\theta}(x^{(i)}) - y^{(i)})$$

~~mx1~~

$$x^{(i)}$$

~~mx(n+1)~~

$$(n+1) \times 1$$

θ

\Rightarrow Vectorized Gradient descent.

$$\begin{aligned} \theta_j &:= \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \\ &\quad \downarrow (n+1) \times 1 \\ &\quad \downarrow (n+1) \times 1 \end{aligned}$$

$$\therefore \theta = \theta - \alpha \frac{1}{m} [(x^{(1)})^T * (h_{\theta}(x^{(1)}) - y^{(1)})]$$

* feature Scaling.

\rightarrow Get every feature into approximately a ~~range~~ $-1 \leq x_i \leq 1$ range.

\rightarrow Examples for Ranges:

maximum range

$-3 \rightarrow 3$

minimum range

$-1/3 \rightarrow 1/3$



~~3) x_1 and x_2 have different ranges.~~ → idea : Make Sure features are on a Similar Scale.

~~2) x_1 and x_2 have different ranges.~~

★ Mean Normalization (Example of feature scaling)

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

(Do not apply to $x_0 = 1$)

where \bar{x}_i = average value of x_i feature

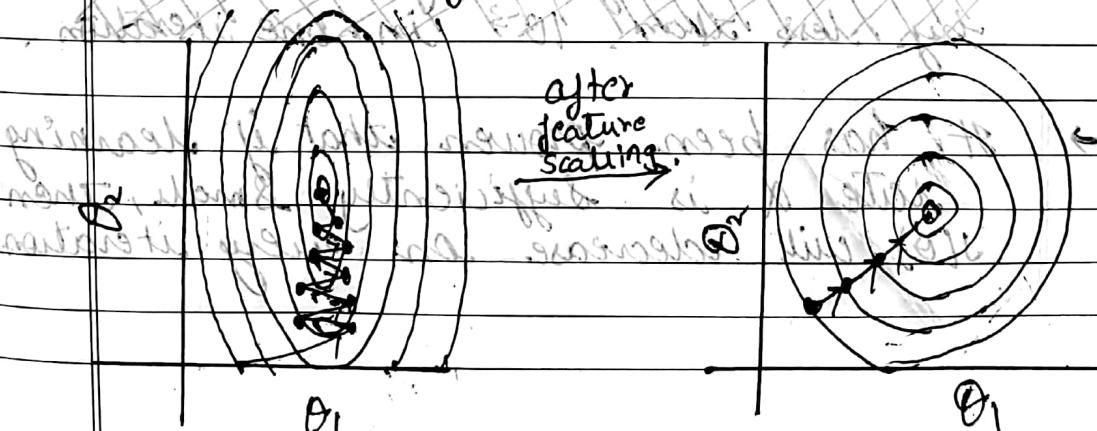
Now,

$$x_i = \frac{x_i - \bar{x}_i}{s_i}, \text{ where } s_i = \frac{\max - \min}{\text{range}}$$

maximum value of feature x_i
minimum value of feature x_i
helps to make range small.
ex: $-3 \leq x \leq 3, -1 \leq x \leq 1, \text{ etc.}$

★ feature scaling on two feature data.

→ See How Cost function look like?



→ Debugging in Gradient descent:

↳ How to make sure gradient descent is working correctly.

(gaining) ↳ How to choose learning rate α .

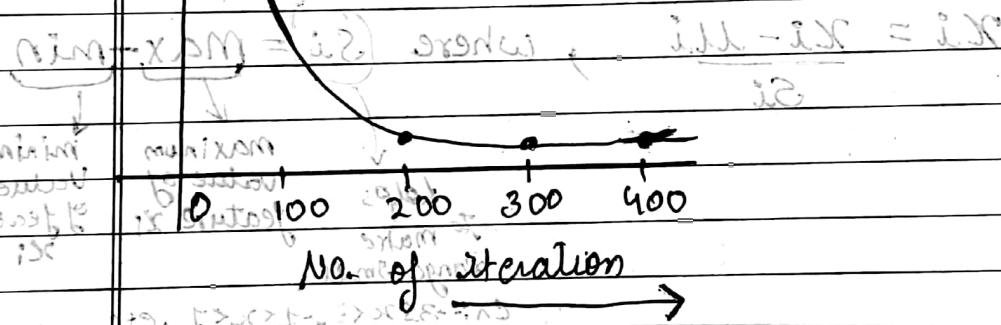
Making sure gradient descent is working correctly

($i=0$) or $\theta_0 = \theta_0 - \alpha J(\theta)$

$$\min J(\theta)$$

start of it is for initial step size = ill consider

with

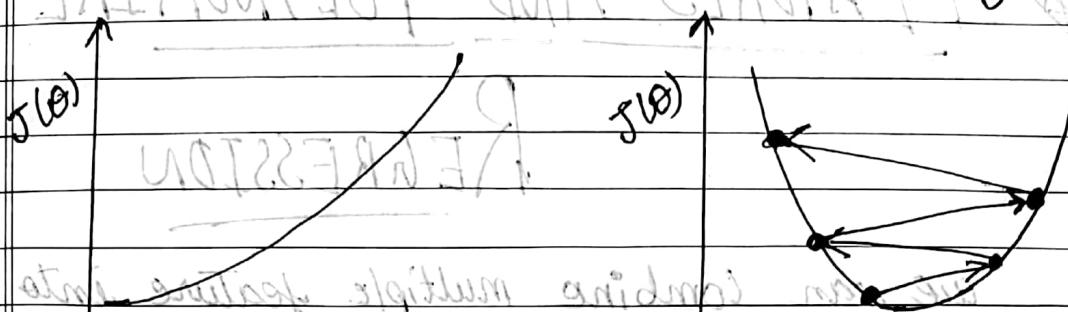


→ $J(\theta)$ should decrease after every iteration

→ Decide convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration.

→ It has been proven that if learning rate α is sufficiently small, then $J(\theta)$ will decrease on every iteration

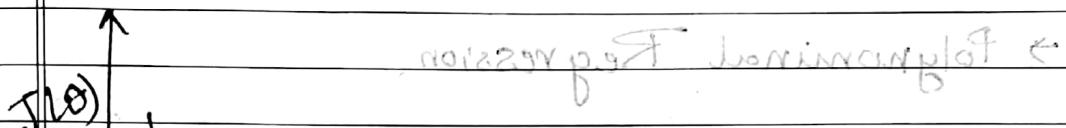
if gradient descent is not working i.e.



are some settings algorithm wouldn't work

either too slow or it stagnates very

so if we want to make it work
we can either increase the step size or decrease the step size



sd. f(x) has narrow parabola zigzagging even
till f'(x) = 0 (minima) so might
have stuck. → sol. try f(x) each
No. of iteration

but it's minimized after 9 epochs. now, we

→ then use a smaller step size. we'll
see much more stable convergence in this

→ But if step size is too small, gradient descent can
be slow to converge.

reduces zigzagging. b. step size is
Summary: α is the step size in
gradient descent. If α is too small,
it will be slow to converge.

- if α is too large $J(\theta)$ may not decrease
i.e. if α is too large, $J(\theta)$ may not converge

To choose α , try

without tool except in the form of
 $\alpha = 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, \dots$

FEATURES AND POLYNOMIAL REGRESSION

we can combine multiple feature into one.
 for example we can combine x_1 and x_2 into a new feature x_3
 by taking $x_1 \cdot x_2$

→ Polynomial Regression.

our hypothesis functions need not be linear (a straight line) if that does not fit the data well.

we can change the behavior or curve of our hypothesis function by making it a quadratic, cubic, or square root function or any other form.

for example, if our hypothesis function is ~~is~~ $h_0(x) = \theta_0 + \theta_1 x$, then we can create addition features based on x_1 to get the quadratic function ($h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$) or the cubic function ($h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$).

To make it a square root function, we could do: $h_0(x) = \sqrt{\theta_0 + \theta_1 x_1 + \theta_2 x_1^2}$

One important thing to keep in mind is,
If your feature is this way then feature scaling becomes very important.

Want to find $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Price (y) vs x_1 (Size (x))

$h_{\theta}(x) = \theta_0 + \theta_1 x_1$

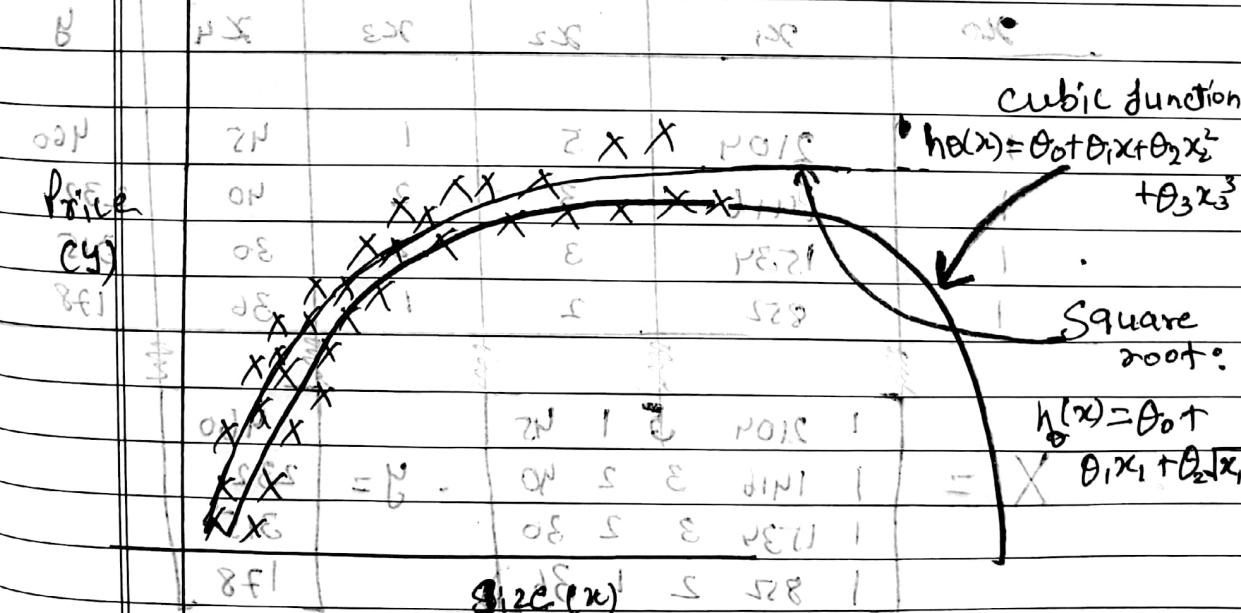
$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2$

Quadratic function

$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3$

Square root:

	High	Medium	Medium	Low	Very Low
(0001)	10000	20000	20000	10000	1000



★ NORMAL EQUATION

Gradient descent gives one way to minimizing $J(\theta)$, let's do this time performing the minimization explicitly and without resorting to an iterative algorithm.

In the "Normal Equation" method, we will minimize $J(\theta)$ by explicitly taking its derivatives with respect to the θ 's and setting them to zero. This allows us to find the optimum theta without iteration. The normal equation formula is given below.

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

Example: m=4

x_0	x_1	x_2	x_3	x_4	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178
$X =$	1 2104 5 1 45				460
	1 1416 3 2 40				232
	1 1534 3 2 30				315
	1 852 2 1 36				178
					$y =$

There is no need to do feature scaling with the normal equation.

GRADIENT DESCENT (GD) V.S. NORMAL EQUATION

→ Need to choose alpha (α)	→ Not Need to choose alpha (α)
→ Needs many iterations	→ Not Need to iterate
→ $O(kn^2)$	→ $O(n^3)$, needs to calculate inverse of $X^T X$
→ Work well when n is large	→ Slow if n is very large

~~if $n > 10,000$? then GD else Normal Equation~~

* What if $X^T X$ is non-invertible?

- Redundant feature (linear dependent)

E.g., $x_1 = \text{size in feet}^2$ $\left. \begin{array}{l} \text{does not} \\ \text{do job} \end{array} \right\} \text{does not} \quad 1m = 3.28 \text{ feet}$
 $x_2 = \text{size in m}^2$ $\left. \begin{array}{l} \text{do job} \\ \text{does not required both features. One can do job.} \end{array} \right\}$

- Too many features (e.g. $m \leq n$)

→ Delete some feature, or use regularization.
Octave: $\text{pinv}(X^T X) * (X^T y)$

XXX Linear Regression: Continuous Value.

Logistic Regression: discrete Value

Date: _____
Page: _____

LOGISTIC REGRESSION

CLASSIFICATION

To attempt classification, one method is to use linear regression and map all prediction greater than 0.5 as a 1 and all less than 0.5 as an 0. However this method does not work well because classification is not actually a linear function.

A non-linear draw is called a sigmoid classification.

Email: Spam / Not Spam

Online Transaction: fraudulent (Yes/ No)?

Tumor: Malignant / Benign?

$y \in \{0, 1\}$ - 0 = "Negative class" (eg; benign tumor)

1 = "Positive class"

(eg., malignant tumor)

XXX y could have more than 2 class
i.e. ($y \in \{0, 1, 2\}$ or $y \in \{0, 1, 2, 3, 4, \dots\}$) etc.

where $y \in \{0, 1\}$ called binary classification.

$$(B \times X) \times (X \times X) \rightarrow \text{Value}$$

* Hypothesis Representation

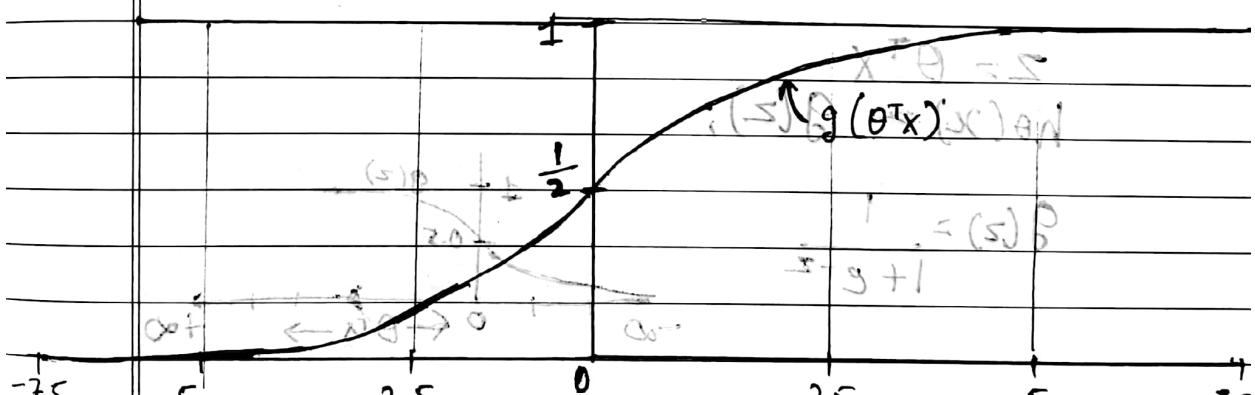
We could approach the classification problem ignoring the fact that y is discrete-valued, and our old linear regression algo. to try to predict y given x . However, it is easy to construct example where this method performs very poorly. Intuitively, it also doesn't make sense for $h_\theta(x)$ to take value larger than 0 or smaller than 0 i.e. $h_\theta(x) > 1$ or < 0 when we know that $y \in \{0, 1\}$. To fix this, one new hypothesis $h_\theta(x)$ to satisfy $0 \leq h_\theta(x) \leq 1$

→ This can be done by plugging $\theta^T x$ into the "logistic function" or "Sigmoid function".

$$h_\theta(x) = g(\theta^T x) \quad \text{--- hypothesis}$$

$$z = \theta^T x$$

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



The function $g(\theta^T x)$ or $g(z)$, $z = \theta^T x$ maps any real number to the $(0, 1)$ interval.

The function $g(\theta^T x)$ or $g(z)$, $z = \theta^T x$ shows that it maps any real number to the $(0, 1)$ interval.

*enlivoo
★ Example to interpret the hypothesis output

$h_{\theta}(x)$ = estimate probability that $y=1$ on input x

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

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

Tell patient that 70% chance of tumor being malignant.

i.e. if probability that it is 1 is 70%, then the probability that it is 0 is 30%

$$P(y=0) + P(y=1) = 1$$

$$\text{or } P(y=0) = 1 - P(y=1)$$

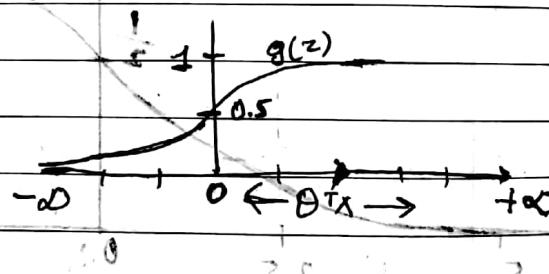


DECISION BOUNDARY

$$z = \theta^T x$$

$$h_{\theta}(x) \approx g(z);$$

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



In order to get our discrete 0 or 1 classification

$$h_{\theta}(x) \geq 0.5 \rightarrow y=1$$

$$h_{\theta}(x) < 0.5 \rightarrow y=0$$

if $z \geq 0$

then $g(z) \geq 0.5$

when Input is greater than or equal to 0. Its output is greater than or equal to 0.5.

Remember

$$z=0 ; c^0 = 1 \Rightarrow g(z) = 1/2$$

$$z=\infty ; c^{-\infty} = 0 \Rightarrow g(z) = 1 + z$$

$$z=-\infty ; c^{(-\infty)} = \infty \Rightarrow g(z) = 0$$

So, if our input to g is $\theta^T x$,

$$h_\theta(x) = g(\theta^T x) \geq 0.5$$

when $\theta^T x \geq 0$

(5) ~~for now we will focus on the first part of the graph~~
for now, these statements can now say;

if $\theta^T x \geq 0 \Rightarrow y = 1$ (if $\theta^T x < 0 \Rightarrow y = 0$)

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

(forward arrow pointing down) : minimum

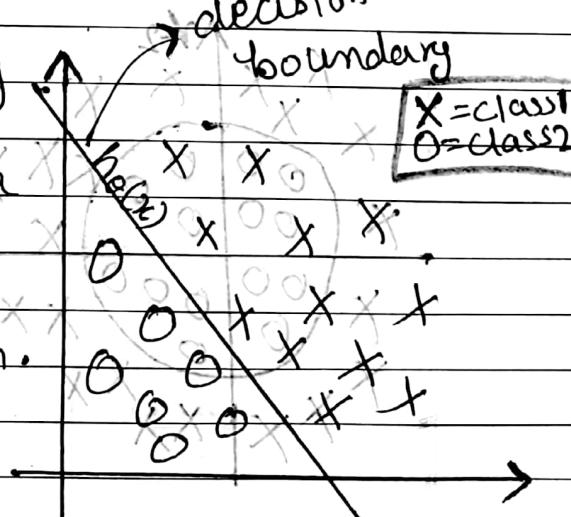
→ The decision Boundary is the line that separates the area where $y=0$ and where $y=1$, it is created by our hypothesis function.

$$0 \leq \theta_0 + \theta_1 x_1 \quad \dots$$

$$1 \leq \theta_0 + \theta_1 x_1 \Leftarrow$$

decision boundary

$X = \text{class 1}$
 $O = \text{class 2}$



Example (logistic function)

$$0 \leq f_i$$

Let, $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}$

for $y=1$ i.e. $\theta^T x \geq 0$

$$\therefore 5 + (-1)x_1 + (0)x_2 \geq 0 \Rightarrow 5 - x_1 \geq 0$$

$$0 = -x_1 \geq -5 \Rightarrow x_1 \leq 5$$

$y=1$	$y=0$
Class 1	Class 2

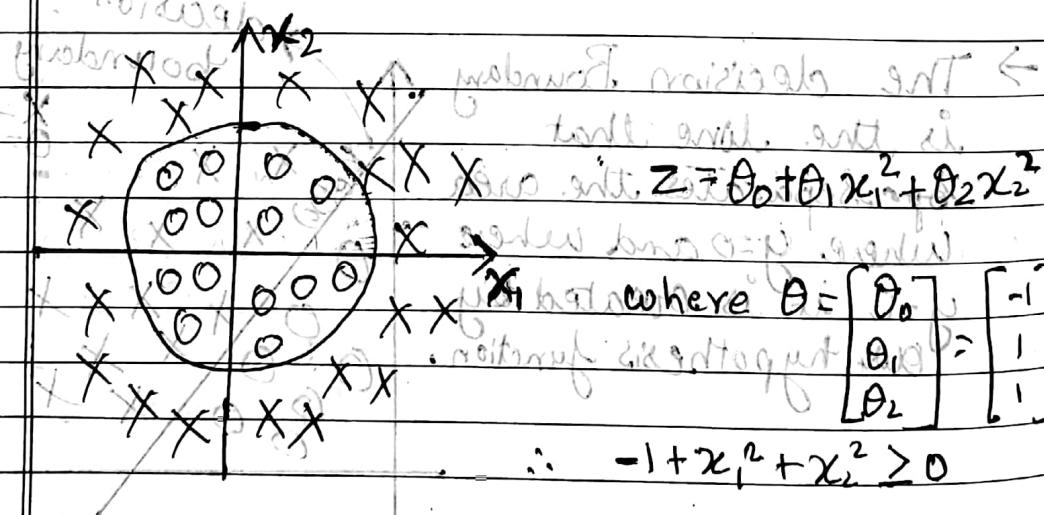
$x^T \theta$ is p of twgn and f_i

$$z.0 \leq (x^T \theta)B = 5 \text{ and } x_1 \rightarrow$$

$$0 \leq x^T \theta \text{ now}$$

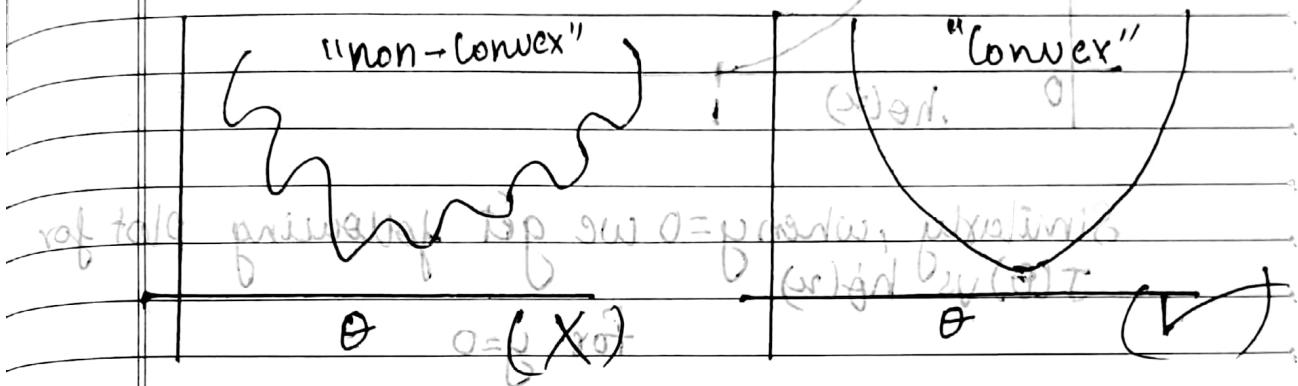
* Again, the input to the Sigmoid function $g(z)$ (e.g. $\theta^T x$) does not need to be linear, and could be function that describes a circle (e.g. $z = \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2$) or any shape to fit data

Example: (non-linear decision boundaries)



Cost FUNCTION

We can't use the same function that we use for linear regression because the logistic function will cause the output to be bumpy, causing many local optima, it will not be convex function.



Logistic regression cost function.

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

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

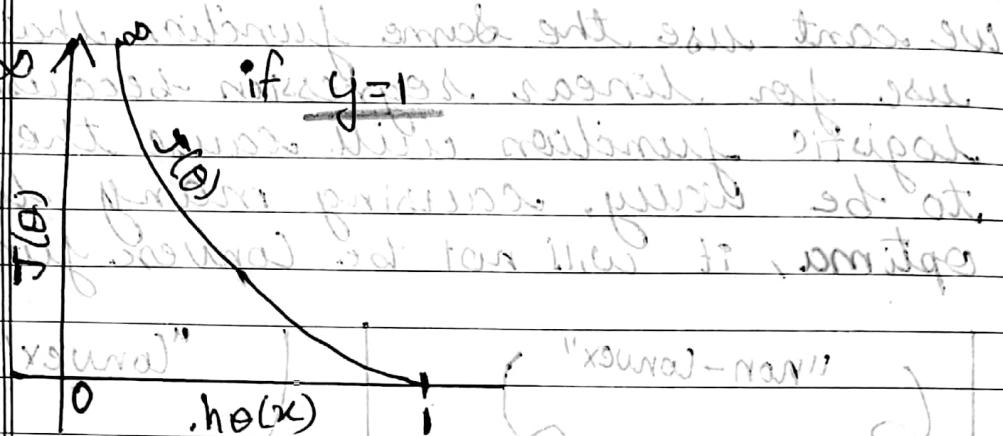
$$\text{Cost}(h_\theta(x^{(i)}), y^{(i)}) = -\log(1-h_\theta(x)) \quad \text{if } y=0$$

Remember: $\infty = (\frac{1}{0} - \log(0))$

$$0 = \log(1-x) \quad \text{if } x = 0 \quad \infty = \frac{1}{1-x} \quad \text{if } x = 1$$

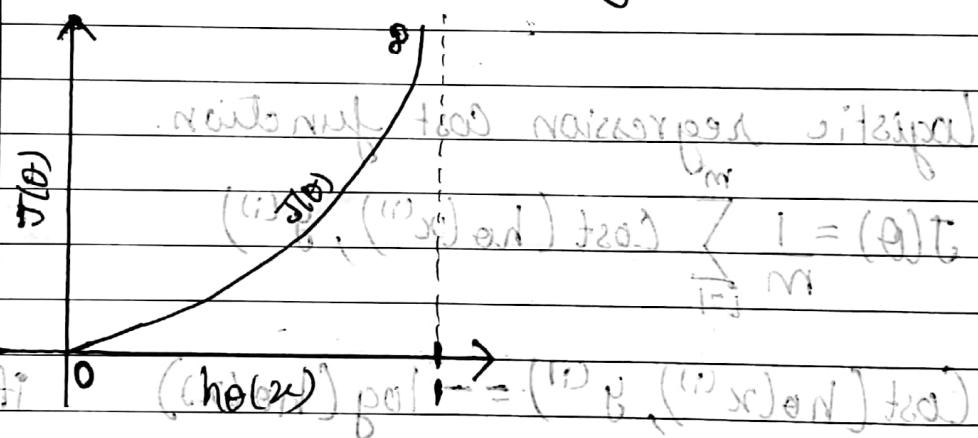
Now just substitute into $J(\theta)$ with both terms if x is 0 or 1.

When $y=1$ we get the following plot for $J(\theta)$ vs $h_{\theta}(x)$



Similarly, when $y=0$ we get following plot for $J(\theta)$ vs $h_{\theta}(x)$

for $y=0$



$$\text{Cost}(h_{\theta}(x), y) = 0 \quad \text{if } h_{\theta}(x) = y$$

$$\text{Cost}(h_{\theta}(x), y) = \infty \quad \text{if } y=0 \text{ and } h_{\theta}(x) \neq 0$$

$$\text{Cost}(h_{\theta}(x), y) = \infty \quad \text{if } y=1 \text{ and } h_{\theta}(x) = 0$$

$(\sigma^{-1})_{pol}$

Note that writing the cost function in this way guarantees that $J(\theta)$ is convex for logistic regression.

* Simplified Cost function.

we can compress our cost functions two conditional cases into one case:

$$(y^{(i)} \ln h_{\theta}(x^{(i)}) + (1-y^{(i)}) \ln (1-h_{\theta}(x^{(i)}))) \quad L = (\theta) T$$

Logistic regression cost function.

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

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

Note: $y=0$ or 1 always

$$\therefore (\theta) T \stackrel{6}{=} \theta - \theta = \theta$$

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

Notice that,

if $y=1$:

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

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

if $y=0$:

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

therefore

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

* Gradient Descent

Goal: Minimize cost function (error) over all data points.

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

$$\left(\text{cost}_j(x^{(i)}) \right) \quad \sum_{i=1}^m = J(\theta)$$

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

$$\left. \begin{aligned} \text{If } \nabla_{\theta} J(\theta) = 0 \\ \theta = \theta \end{aligned} \right\} = (\nabla_{\theta} J(\theta))$$

Algorithm: $\theta^{(k+1)} = \theta^{(k)} - \alpha \nabla_{\theta} J(\theta)$

Repeat until $\nabla_{\theta} J(\theta) = 0$: Stop

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

$$\left. \begin{aligned} (\theta^{(k+1)})_j &= (\theta^{(k)})_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \\ &= (\theta^{(k)})_j - \alpha \nabla_{\theta_j} J(\theta) \end{aligned} \right\}$$

3

(Simultaneously updating all θ_j).

here,

$$\left. \begin{aligned} (\theta^{(k+1)})_j &= \frac{\partial}{\partial \theta_j} J(\theta) \\ &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \end{aligned} \right\}$$

$$\left. \begin{aligned} (\theta^{(k+1)})_j &= (\theta^{(k)})_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \\ &= (\theta^{(k)})_j - \alpha \nabla_{\theta_j} J(\theta) \end{aligned} \right\}$$

∴ Repeat L

$$\left. \begin{aligned} \theta_j &:= \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ &= \left((\theta^{(k+1)})_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \left((\theta^{(k)})_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) \end{aligned} \right\}$$

Notice that this algorithm is identical to the one we used in linear regression.

algorithm is identical, but definition is different (partial derivative of linear regression cost function and partial derivative of logistic regression cost function is totally different).

* Vectorized implementation of:

1) Cost function.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[-y_i \log(h_i) - (1-y_i) \log(1-h_i) \right]$$

[2] Gradient descent ...

$$\theta := \theta - \alpha \sum_{i=1}^m x^T (g(x\theta) - \vec{y})$$

* ADVANCED OPTIMIZATION

→ Gradient descent

→ Conjugate gradient → are more sophisticated, faster ways to optimize

→ BFGS → that can be used instead of gradient

→ L-BFGS, Nelder-Mead → descent, "Conjugate" or "BFGS" But are more complex

This page shows how to use
in-build Advanced optimization function
on octave. ("fminunc()")

*enlive

Date: _____
Page: _____

Octave provides highly optimized library.
("fminunc()")

We first need to provide a function the
evaluates the following two function
for a given input value θ .

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

We can write a single function that return
both of these:

function [JVal, gradient] = CostFunction(theta)

JVal = [... code to compute $J(\theta)$];
gradient = [... code to compute $\frac{\partial}{\partial \theta} J(\theta)$]

end

Now, we can use "fminunc()" optimization algo.
along with "optimset()" junction that creates an
object containing the options we want to send
to "fminunc()".

options = optimset('GradObj', 'on', 'MaxIter', 100);
initialTheta = zeros(2, 1);

[optTheta, functionVal, exitFlag] = fminunc(@CostFunction, initialTheta, options);

we give our cost function, our initial vector of theta
and "options" object to "fminunc()".

* MULTI-CLASS CLASSIFICATION:

$(\theta; x) \hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$

ONE-VS-ALL

$(\theta; x) \hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$

* multi-class classification.

→ Email foldering / Tagging : work / friends / family / Hobby.

$y \in \{0, 1, 2, 3\}$

→ Medical diagnosis: Normal, cold, Flu
 $y \in \{0, 1, 2\}$.

→ Weather: Sunny, cloudy, Rain, Snow

with weight $y \in \{0, 1, 2, 3\}$. Signal ment.

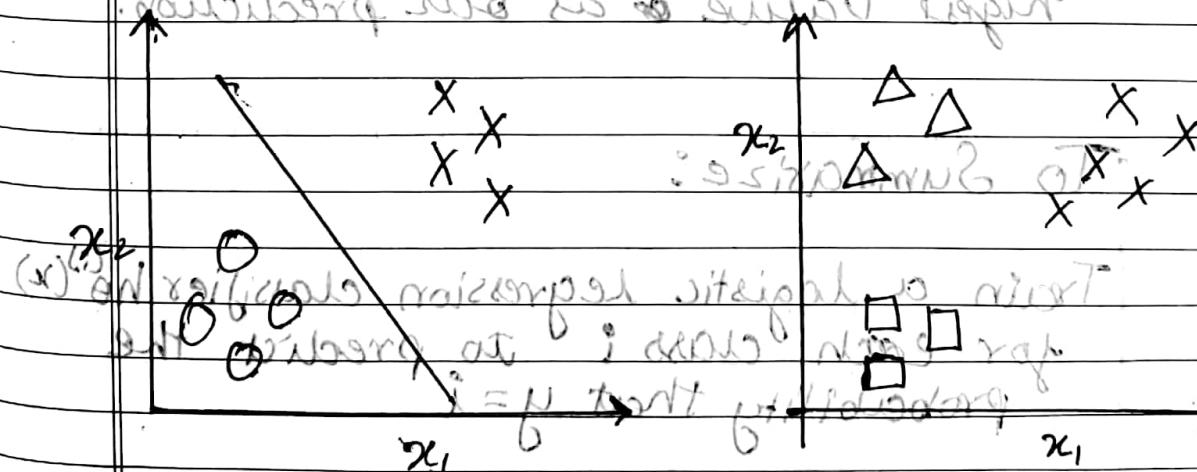
city ab. sun. evn. rain & snow

situation: pressure, wind, etc.

Binary classification vs. multi-class

either true / false or classification

multiple classes as multi-label



$(x)_i$ on x_{ij}

Let $y \in \{0, 1, \dots, n\}$

$$h_{\theta}^{(0)}(x) = P(y=0 | x; \theta)$$

$$h_{\theta}^{(1)}(x) = P(y=1 | x; \theta)$$

$$h_{\theta}^{(2)}(x) = P(y=2 | x; \theta)$$

$$h_{\theta}^{(n)}(x) = P(y=n | x; \theta)$$

$$\text{prediction} = \max(h_{\theta}^{(0)}(x), h_{\theta}^{(1)}(x), \dots, h_{\theta}^{(n)}(x))$$

We are basically choosing one class and then lumping all the others into single second class. Due to this repeatedly applying binary logistic regression to each case, and then use the hypothesis that return the highest value as our prediction.

To Summarize:

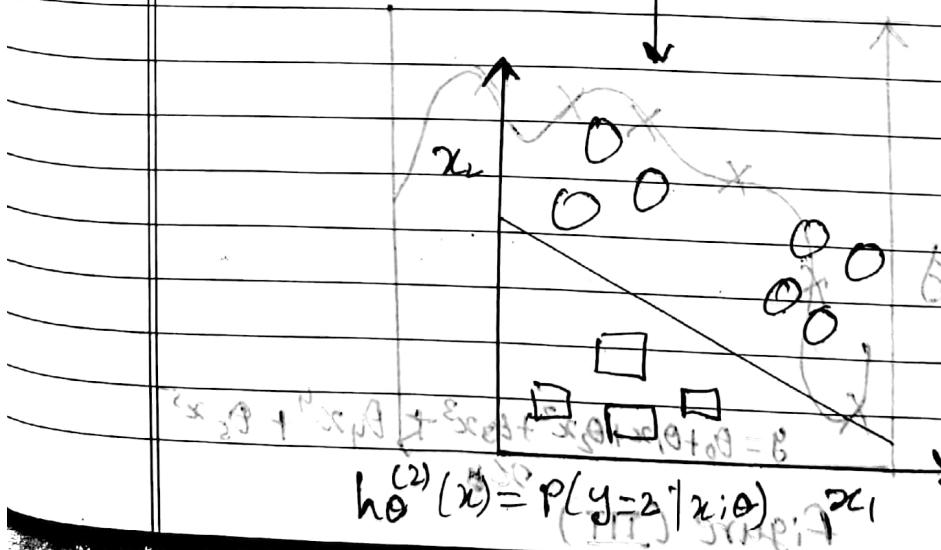
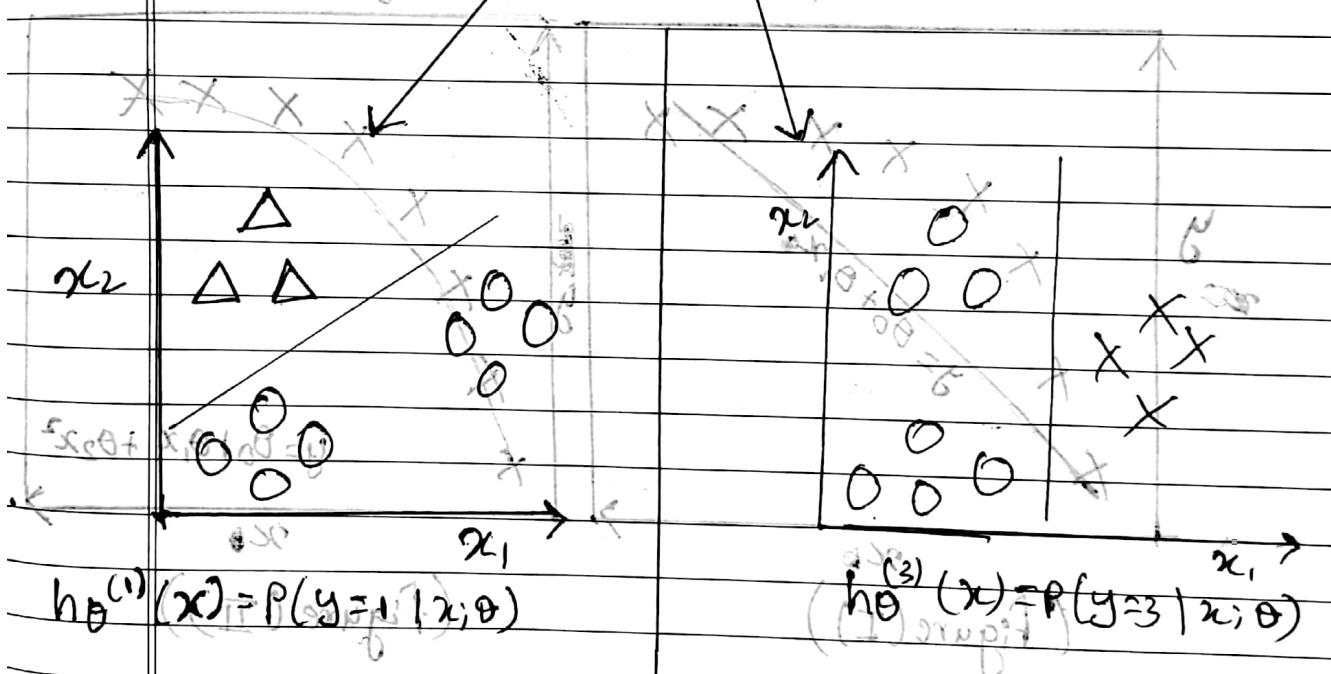
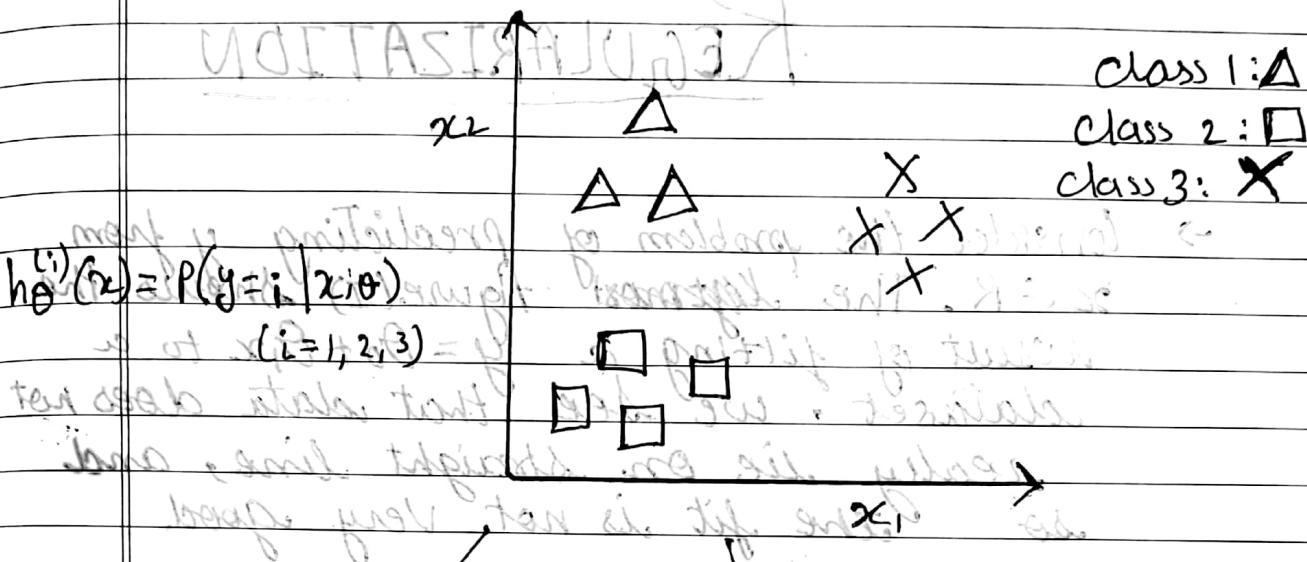
Train a logistic regression classifier $h_{\theta}^{(i)}(x)$ for each class i to predict the probability that $y=i$.

On a new input x to make a prediction, pick the class i that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

Example:

One-vs-all (One-vs-rest)

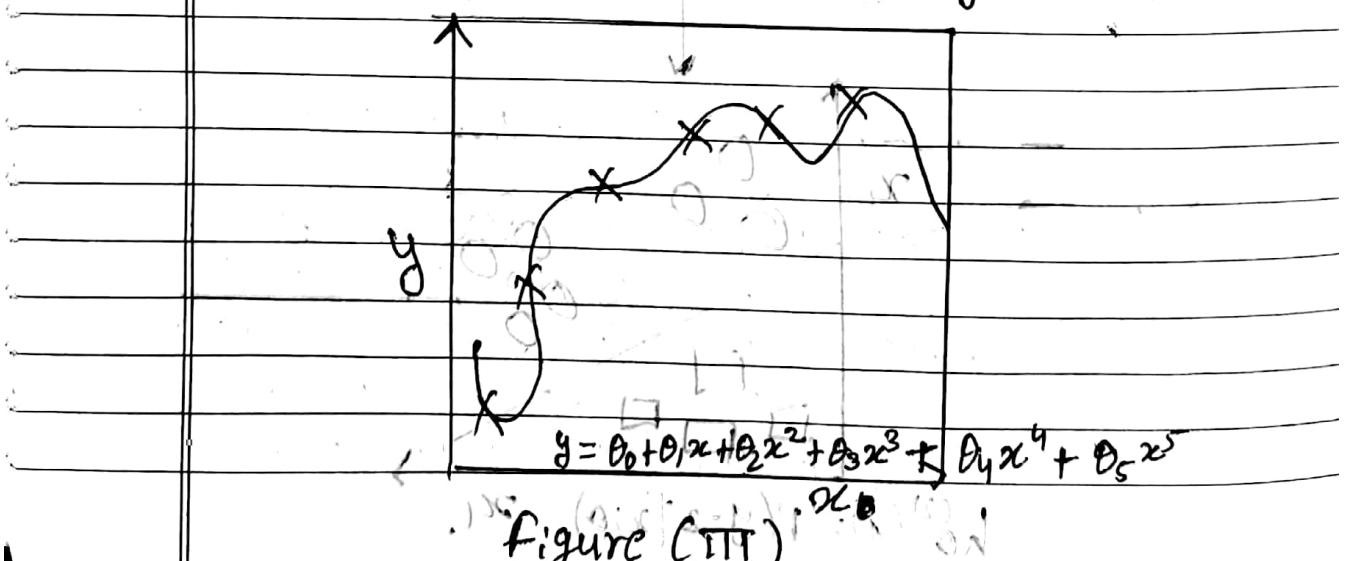
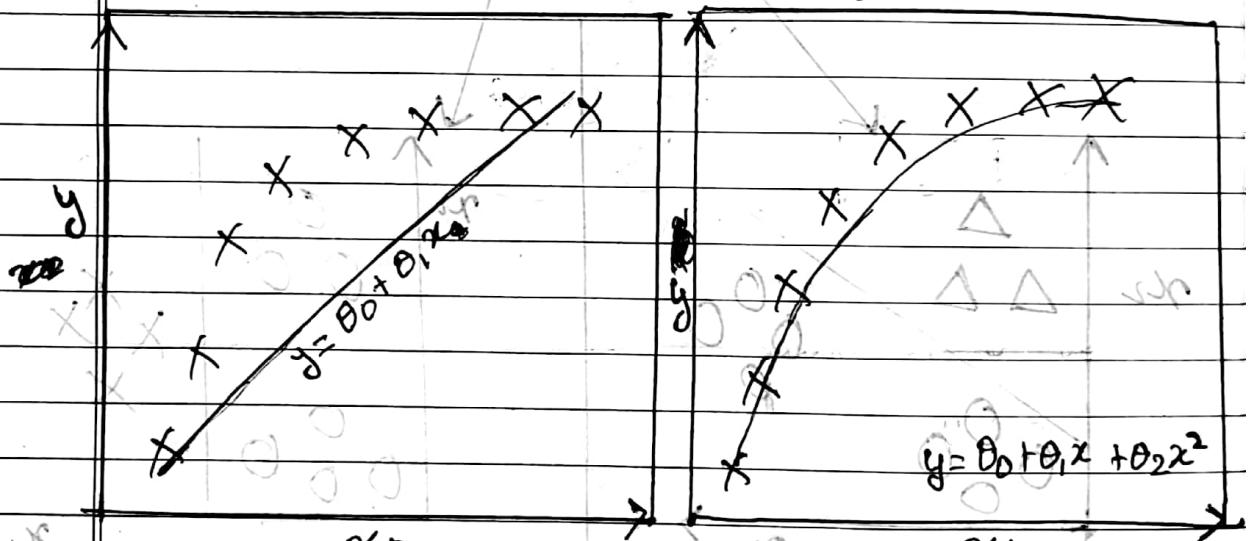


* THE PROBLEM OF OVERFITTING :

(too many terms) $y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_n x^n$

REGULARIZATION

→ Consider the problem of predicting y from $x \in \mathbb{R}$. The figure(I) shows the result of fitting a $y = \theta_0 + \theta_1 x$ to a dataset. We see that data does not really lie on straight line, and so the fit is not very good.



instead, if we had added an extra feature x^2 , and fit $y = \theta_0 + \theta_1 x + \theta_2 x^2$, then we obtain a slightly better fit to the data (figure (II)). Naively, it might seem that the more features we add, the better. However, there is also a danger in adding too many features. The figure (figure (III)) is result of fitting a 5th order polynomial $y = \sum_{j=0}^5 \theta_j x^j$. We see that even though the fitted curve passes through the data perfectly we would not expect this to be a very good predictor of, say (y) for different (x) . Without formally defining what these terms mean, we'll say the figure (I) shows an instance of Underfitting — in which the data clearly shows a structure not captured by the model — and the figure (III) is an example of Overfitting.

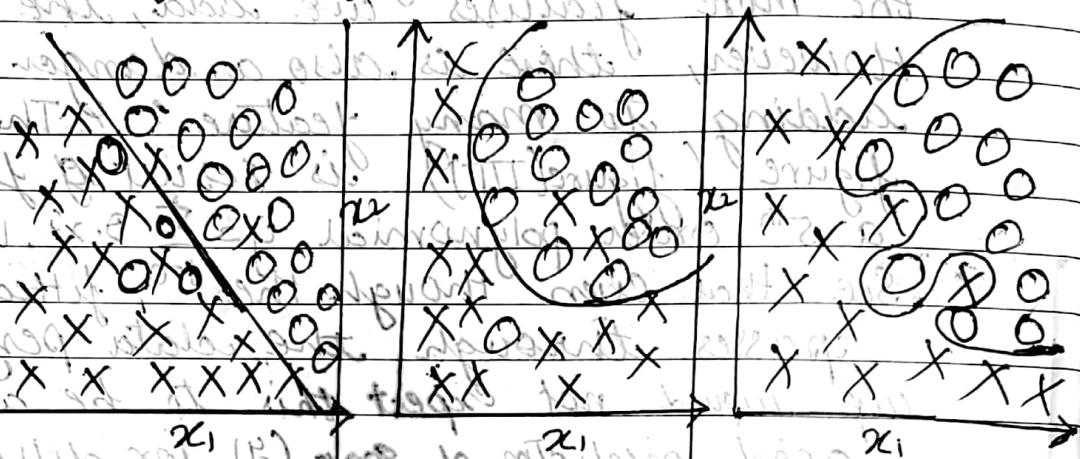
Bias Variance

Underfitting, or high bias, is when the form of our hypothesis function h maps poorly to the trend of the data. It is usually caused by a function that is too simple or uses too few features.

Overfitting or high variance, is caused by a hypothesis function that fits the available data does not generalize well to predict new data. It is usually caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.

This terminology is applied to both linear and logistic regression.

* Logistic regression:



$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 + \theta_6 x_1^3 x_2 + \dots)$$

\rightarrow Underfitting

\rightarrow Just right

\rightarrow Overfit.

* Addressing Overfitting.

1. Reduce number of features.

- manually select which features to keep.

- Model Selection algorithm

2. Regularization.

Keep all the feature, but reduce magnitude / New values of parameters θ_j

- Works well when we have a lot of features, each of which contributes a bit to predicting y .

* Cost function

If we have overfitting from our hypothesis function, we can reduce the weight that some of the terms in our function carry by increasing their cost.

Say we wanted to make the following function more quadratic.

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

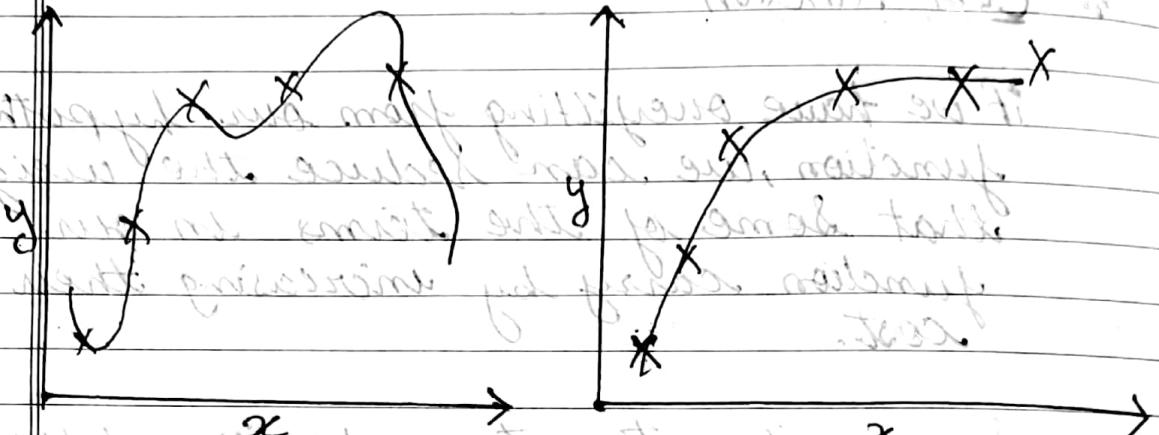
We'll want to eliminate the influence of $\theta_3 x^3$ and $\theta_4 x^4$. Without actually getting rid of these features or changing the form of our hypothesis, we can instead modify our cost function.

$$\text{min}_{\theta} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \cdot \theta_3^2 + 1000 \cdot \theta_4^2.$$

We've added two extra terms at the end to

inflate the cost of θ_3 and θ_4 . Now, in order for a cost function to get close to zero, we will have to reduce the value of θ_3 and θ_4 to near zero. This will in turn greatly reduce the value of $\theta_3 x^3$ and $\theta_4 x^4$ in our hypothesis function.

As a result, we see that new hypothesis looks like a quadratic function but fits the data better due to the extra small terms $\theta_3 x^3$ and $\theta_4 x^4$.

~~Intuition~~

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$\theta_3 \approx 0 \quad \theta_4 \approx 0$$

we penalize and make θ_3, θ_4 really small.

$$\min_{\theta} \sum_{i=1}^{m/2} (\hat{y}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2$$

→ above ~~Intuition~~ Intuition Shows how we can penalize and make value of θ very small, ≈ 0 , in order to fit ~~hypothesis~~ hypothesis on data ~~perfect~~ just right.

We could also regularize all of our theta parameters in a single summation as:

$$\min_{\theta} \sum_{i=1}^{m/2} (\hat{y}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$$

The λ , or lambda, is the regularization parameter, it determines how much the costs of our theta parameters are inflated.

- # Small value for parameters $\theta_0, \theta_1, \dots, \theta_n$
- "Simpler" hypothesis
 - Less likely to overfitting.

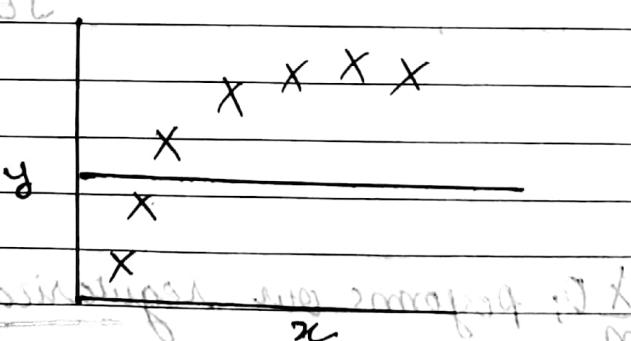
Using the ~~as~~ this cost function with extra summation, i.e. can smooth the output of our hypothesis function to reduce overfitting. If λ is chosen to be too large, it may smooth out the function too much and cause underfitting.

In regularized linear regression, we choose θ to minimize:

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (\theta_0 x^{(i)}) - y^{(i)} \right]^2 + \lambda \sum_{j=1}^n \theta_j^2$$

If λ is set to an extremely large value, (say $\lambda = 10^{10}$)

- Algorithm work fine; setting λ to be very large ^{can't hurt}
- Algorithm fails to eliminate overfitting.
- Algorithm results in underfitting. (fails to fit even from training data well).
- Gradient descent will fail to converge.



~~except~~ θ_0 will be $\theta_0 \approx 0$, ($j = 1, 2, \dots, n$)
except θ_0 , and
~~cause~~ underfit.
~~cause~~ underfit.

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$\theta_1 \approx 0$ $\theta_2 \approx 0$ $\theta_3 \approx 0$ $\theta_4 \approx 0$

Regularized Linear Regression

→ Cost function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\text{h}\phi(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$$

$\min_{\theta} J(\theta)$ where θ_0 is excluded (The bias term)

→ Gradient Descent

we will modify our gradient descent function to separate out θ_0 from the rest of the parameters because we do not want to penalize θ_0 .

Repeat i

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

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (\text{h}\phi(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \right] + \frac{\lambda \theta_j}{m}$$

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

The term $\frac{\lambda \theta_j}{m}$ performs our regularization.