

Chapter 1 :

Machine learning:

Field of study that gives computers the ability to learn without being explicitly programmed

-Arthur Samuel

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

-Tom Mitchell.

Machine learning is an application of Artificial intelligence (AI) that provides the systems the ability to automatically learn and improve from experience without being explicitly programmed.

Machine learning allows software applications to become more accurate in predicting outcomes without being explicitly programmed.

Example

definition

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 .

↳ LP-6

Suppose you email watches which email you do or do not mark as spam and based on that learns how to better filter spam.

- Classifying emails as spam or not spam - Task (T)
- Labelling you label emails as spam or not - Exp (E)
- The no. of emails correctly classified. - Performance (P)

2) A robot driving learning problem:

T : Driving on public four lane highway using vision sensors

P : Avg. distance traveled before an error (as judged by human overseer)

E : A sequence of images and steering commands recorded while observing a human driver

3) A checkers learning Problem:

T : playing checkers

P : percent of games won against opponents

E : Playing practice games against itself.

Ex:

Prob

Prob

Machine learning Algorithms

→ Supervised Learning

→ Unsupervised Learning

Supervised learning:

In machine learning it's a task of inferring function from labeled training data. Training data consist of a set of training examples.

Here, each example is a pair consisting of a input object (vector) and a desired output.

A supervised learning algorithm analyzes the training data and produces an inferred function which can be used to map new examples likewise.

Supervised learning algorithm produces more right answers and often referred as "Regression" problem

Regression: Predict continuous valued output

Supervised learning → Regression

→ classification

Ex: You're running a company and you want to develop learning algorithms to address each of two problems

Problem 1: You have a large inventory of identical items. You want to predict how many of these items will sell over next 5 months.

Problem 2: You'd like software to examine individual customer accounts, and for each account decide if it has been hacked/compromised.

Problem 1 is a regression problem

bcz suppose we have 1000's of value & u can treat it as real/continuous value

Problem 2 is a classification problem

set values 0 - Not hacked

1 - hacked. and predict the algorithm

→ Approach: steps to solve a supervised learning problem

1. Determine the type of training examples
2. Gather a training set. ie gathering a set of input objects and corresponding output
3. Determine the input feature representation of the learned function
4. Determine the structure/ flow of the learning algorithm
5. Complete the design. ie run the learning algorithm on the gathered training set
6. Evaluate the accuracy of the learned function.

Unsupervised learning:

It's a machine learning task of inferring a function to describe unlabeled data (classification or categorization is not included).

Since the ex. given are unlabeled there is no evaluation of accuracy of the structure.

Approaches to solve unsupervised learning may be clustering, Neural networks, etc.

Examples

- A highly developed AI that serves as a house keeping robot develops a theory that there is usually dust under sofa. Each week, the theory is confirmed as the robot often finds dust under the sofa. Nobody explicitly tell the robot the theory is correct but it develops on its own.
- Given a set of news article found on the web, group them into set of articles of the same story
- Given a database of customer data, automatically discover market segments and group customers into different market segments.

2. Vic

Semi Supervised learning:

- It's a class of supervised learning. It lies between supervised and unsupervised learning.
- It has small amount of labeled data and a huge amount of unlabeled data.
 - Labeled and unlabeled data when used in conjunction with labeled data, considerably improves learning accuracy.

possible
They can

Reinforcement learning:

This differs from standard supervised learning in which correct input/output pair are never present, nor sub optimal action explicitly corrected. Instead the focus is on performance which involves finding balance betⁿ parameters.

Reinforcement learning allows machines & software agents to automatically determine the ideal behaviour within a specific context in order to maximize its performance.

IP-2

Application of Machine learning

1. Virtual Personal Assistants

Siri, Google, Cortona are some popular examples. As the name suggests they assist in finding information when asked over voice.

For answering the personal assistant look out for the information stored or recalls your related queries or send a command to other resources to collect info.

ML is an important part of these personal assistants as they collect & refine the info on the basis of your previous involvement with them.

3. S

lea
A

2. Video Surveillance

These systems now a day make it possible to detect crime before they happen. They track unusual behaviour of people. The sys can thus gives an alert to human attendants.

3. Social Media Services.

Social Media platforms utilize machine learning for their own and user benefits.

A few example include:

→ People You May Know: ML works on simple concept understanding with experience. Facebook continuously notices the friend you connect, the often profile you visit, your interest work place or group that you share, and on this basis suggest frndz

→ Face Recognition: When uploaded a picture with a friend, the Facebook instantly recognizes that friend and check then match with the people in your friend list

→ Similar Pins: ML is the core element of Computer vision. Pinterest uses computer vision to identify the objects or pins in the image and recommend similar pins

4. Email Spam & Malware filtering

5. Search Engine Result Refining

6. Product Recommendations

Chapter - 2 Supervised Learning

* Linear Regression = linear Regression.

In supervised learning there is basically a data set where 'x' is the input and 'y's are the corresponding output, 'm' is the no. of training examples, and the job is to predict the output which is also known as hypothesis function

(x, y) - denotes single training example

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

Ex: Size in feet² (x)

2104
1416
1534
852
:

Price in 1000's \$ (y)

460

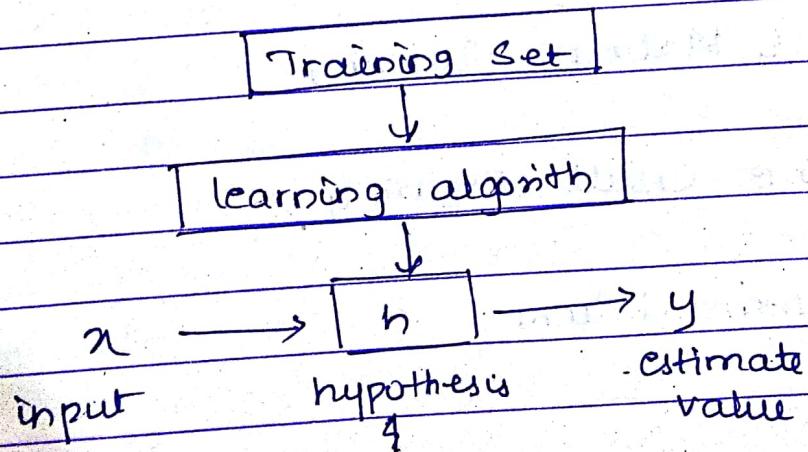
232

815

178

$$x^1 = 2104 \text{ with } y^1 = 460$$

working of supervised learning algorithm



h - function that maps x 's to y 's

Hypothesis \hat{y}

Representing hypothesis :

It's an estimate or model or function with limited inputs/ training examples / features to which an inference is drawn at that instant.

It's a proposed explanation made on the basis of limited evidence as a starting point for further investigation

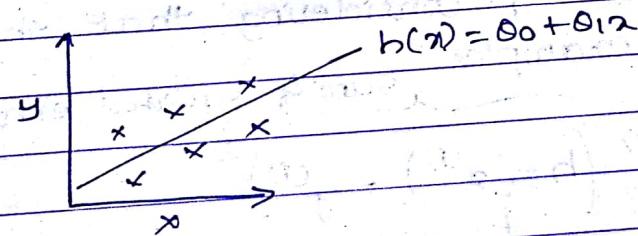
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$\theta_0, \theta_1 \rightarrow$ parameters/weight

$h_{\theta}(x) \rightarrow$ predicted o/p.

$b(x)$

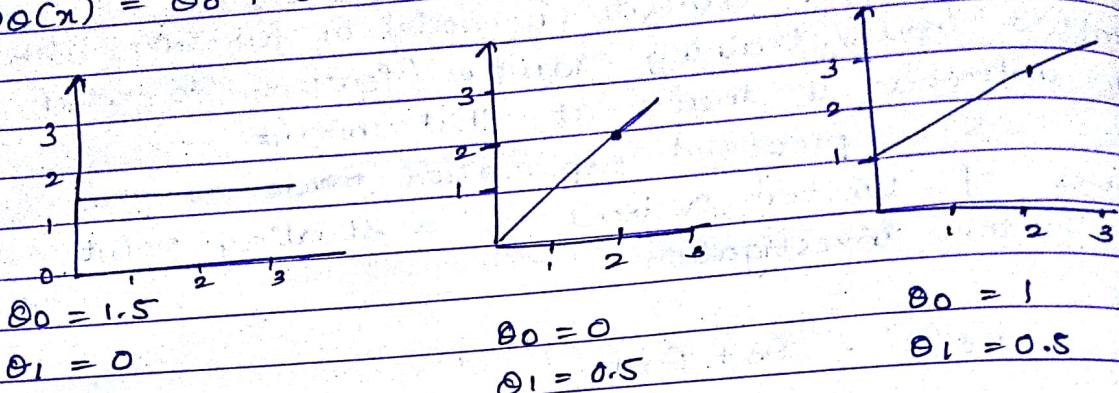
$h_{\theta}(x)$ means that we are predicting y some straight line function. straight or linear function



So this regression is having only one variable x there is no i 'th index as such hence the name linear regression with one variable or also known as Univariate linear regression

Cost functions

$$h_0(x) = \theta_0 + \theta_1 x$$



The update



θ_0 and θ_1 values
are chosen so that
 $h_0(x)$ is close to y for
training ex (x, y)

minimizing θ_0 and θ_1 , considering that there
are m training examples.

$$\text{minimize}_{\theta_0, \theta_1} = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

↓
avg

$\underbrace{\theta_0 + \theta_1 x^{(i)}}_{\text{sum of squared errors}}$

rewriting

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

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

$\underbrace{\quad}_{\text{cost function}}$

cost fun
function

The gr
when θ_0 is

→ Understa

$$h_0(x) =$$

J(

fun

pa

→

A

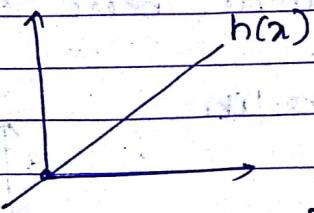
cost function is also called as "squared error function"

The goal is to minimize the cost function is $J(\theta_0, \theta_1)$ where θ_0 is the bias function.

→ Understanding Cost function using simplified hypothesis function

$$h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \text{function of } x$$

$$\theta_0 = 0$$



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

function of parameter θ_1

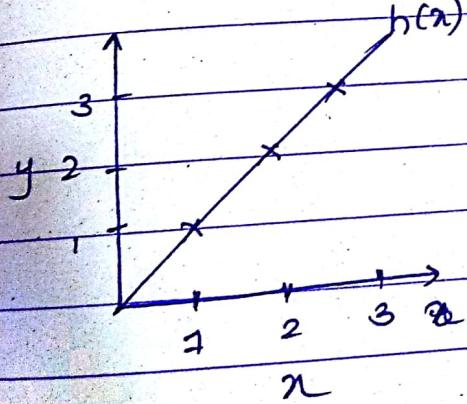
→ $h_{\theta}(x)$

$$\text{Assume } \theta_0(x) = 1$$

$$\theta_1 = 1$$

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

Waiting cost function



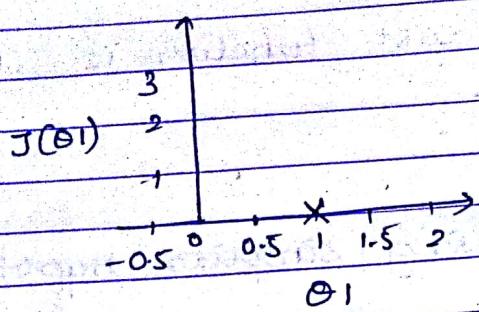
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

$$= \frac{1}{2m} (0^2 + 0^2 + 0^2)$$

$$J(1) = 0$$

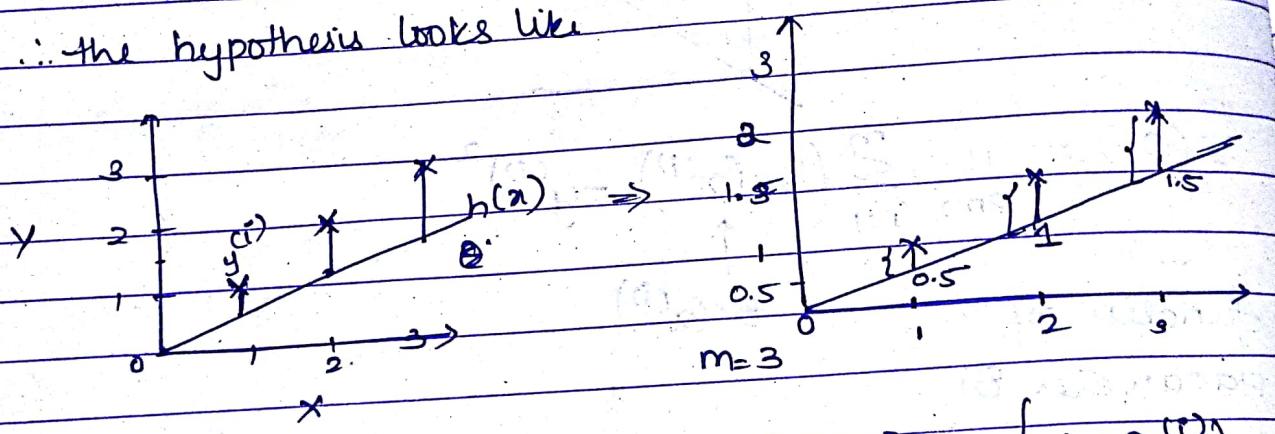
plotting $J(\theta_1)$



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

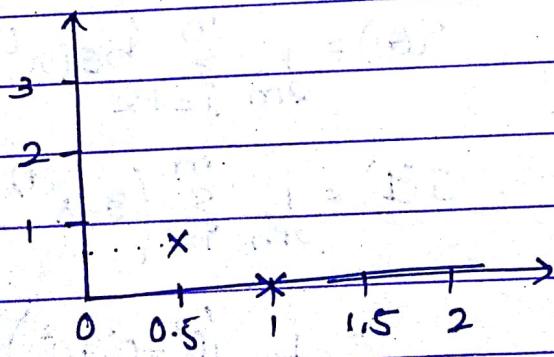
consider another example where $\theta_1 = 0.5$

\therefore the hypothesis looks like

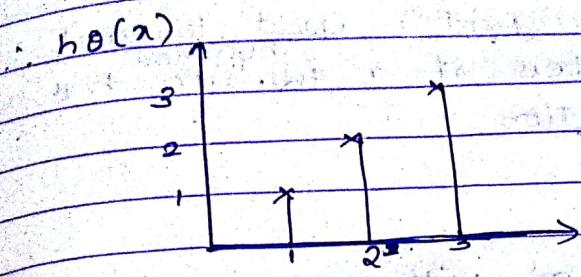


$$J(0.5) = \frac{1}{2m} \left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right]$$
$$= \frac{1}{6} [3.5]$$

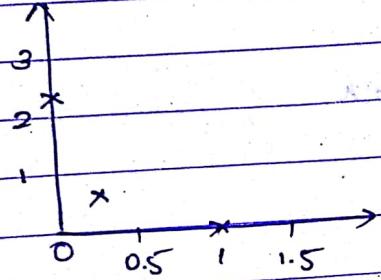
$$J(0.5) \approx 0.58$$



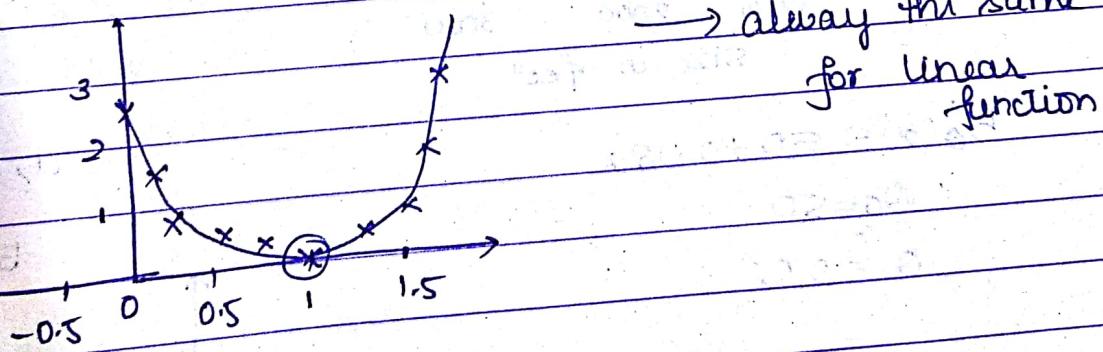
consider $\theta_1 = 0$



$$\begin{aligned} J(0) &= \frac{1}{2m} [1^2 + 2^2 + 3^2] \\ &= \frac{1}{6} [14] \\ &= 2.3 \end{aligned}$$



the $J(\theta)$ can be computed for different values
and that now will look like



Since the goal was to minimize the $J(\theta_1)$ the cost function
is minimum at $\theta_1 = 1$

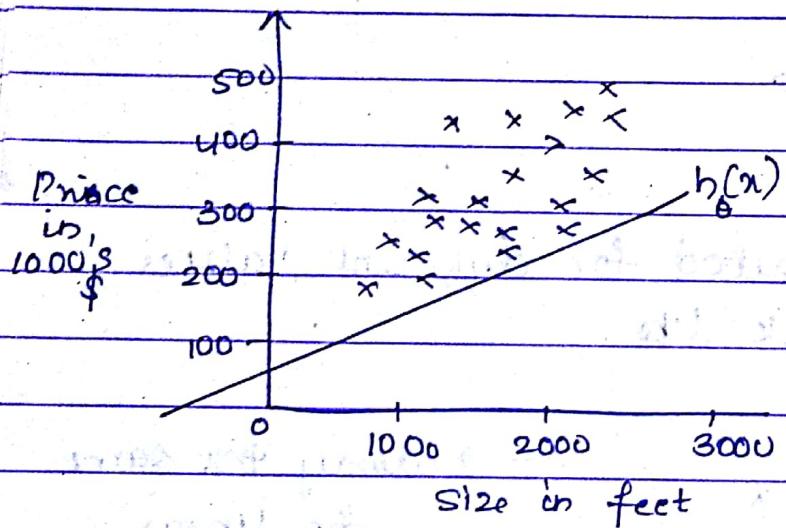
Gradient descent

Understanding cost function with cost function having more than 1 θ.

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

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

considering a training ex.

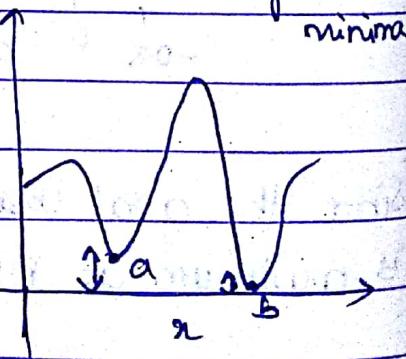


$$h_{\theta}(x) = 50 + 0.06x$$

$$\theta_0 = 50$$

$$\theta_1 = 0.06$$

plotting corresponding $J(\theta_0, \theta_1)$



distance betⁿ n and point illj with point b gives the $J(\theta_0)$ and $J(\theta_1)$ respectively

Gradient Descent

It's an optimization algorithm used to find the values of parameters of a function/hypothesis that minimizes the cost function

$J(\theta_0, \theta_1) \rightarrow$ given function

$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ - required.

\therefore start with some θ_0, θ_1

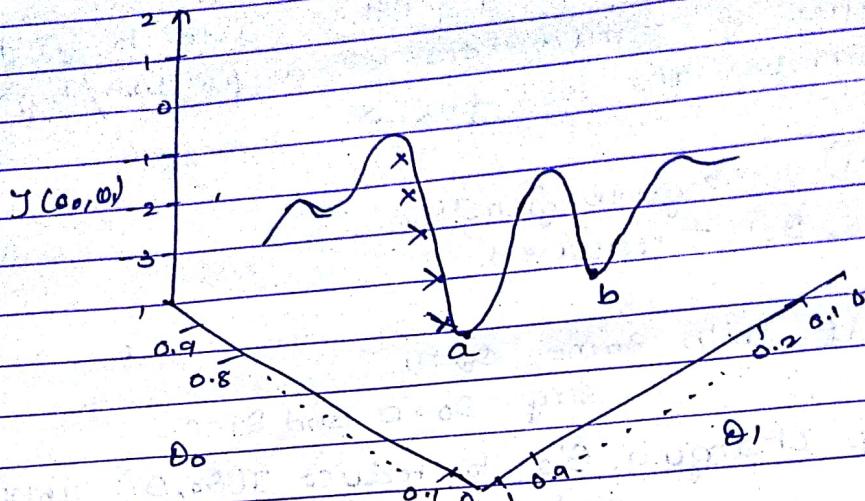
say $\theta_0 = 0$ and $\theta_1 = 0$

and keep changing θ 's to reduce $J(\theta_0, \theta_1)$ until you end up at a minimum

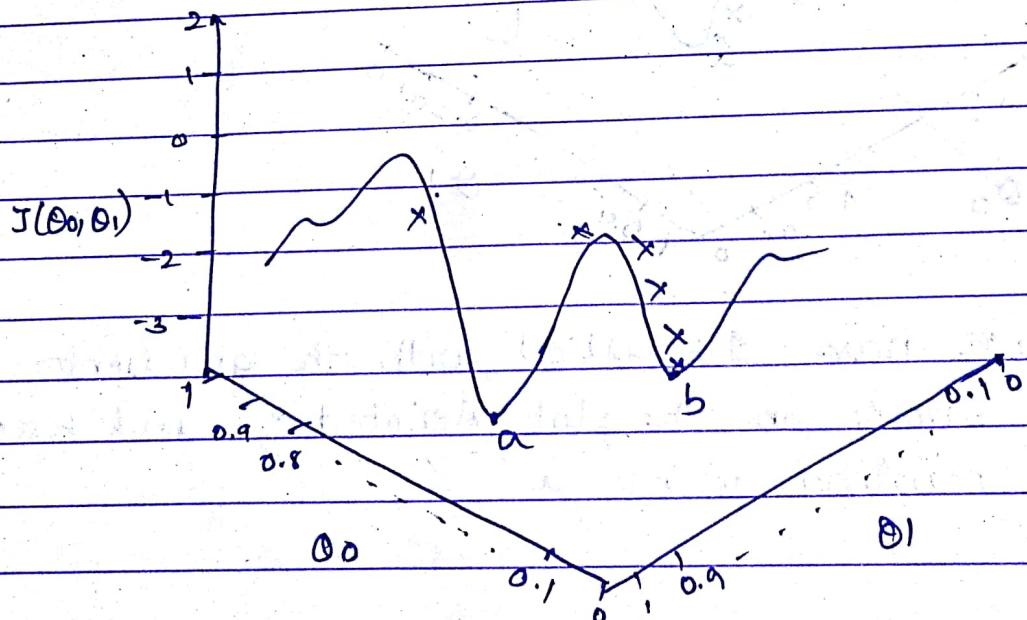
Gradic

Repe

sim



Imagine we have started with a gradient descent as shown in the above plot (denoted x) and reach a local minima point 'a'



Consider another gradient descent slightly away from the previous one and taking certain steps ends up reaching at another local minima point 'b'

Gradient descent algorithm

{ := Assignment

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \left\{ \begin{array}{l} j=0 \text{ and } j=1 \\ \end{array} \right.$$

}

Simultaneous update (Correct Method)

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

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

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

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

α - learning rate / decides step size

α is large the ~~or~~ diff betⁿ two θ values
is more larger steps & vice versa

~~θ₀~~
~~θ₁~~

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

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

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

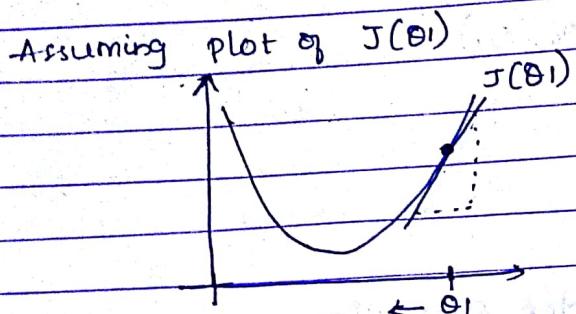
$$\theta_1 := \text{temp}_1$$

* Action of derivative term on θ considering the gradient descent algorithm.

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \left\{ \begin{array}{l} j=0 \text{ & } j=1 \\ \text{derivative} \end{array} \right.$$

To understand the derivative considering a simple example

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



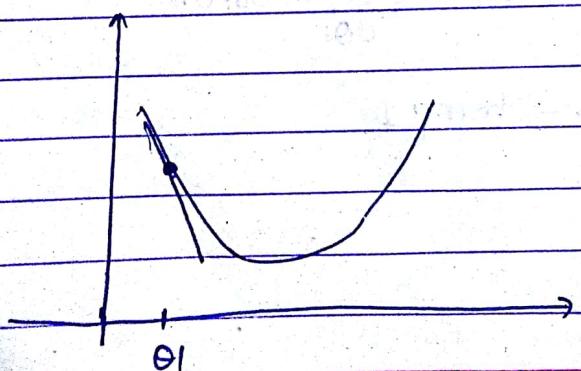
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

The derivative indicates
the tangent drawn to the
function $J(\theta_1)$ function is a slope
ii) the line drawn has the slope
ie $\frac{\partial}{\partial \theta_1} J(\theta) \geq 0$

: the update θ value will be
 $\theta_1 := \theta_1 - \alpha (\text{+ve number})$

here the θ values are decreased to reach a local minima

* Considering another θ_1 for the same function $J(\theta_1)$



the derivative term
changes ~~as~~.

the tangent drawn
is -ve slope
 $\text{ie } \frac{\partial}{\partial \theta_1} J(\theta) \leq 0$

the update θ value will be

$$\theta_1 := \theta_1 - \alpha (\text{-ve number})$$

The θ value is increased to reach a local minima

Action of Le

i) α is

$\theta_1 = \theta$

If the

by a very

near to

This req

minimum

is reached

the update

is zero

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

the update

is zero

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

the update

is zero

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

the update

is zero

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

the update

is zero

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

the update

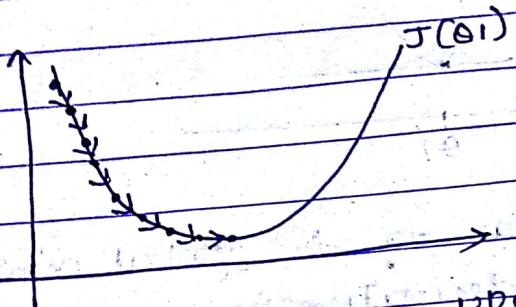
is zero

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

Actions of Learning Rate (α) on the θ .

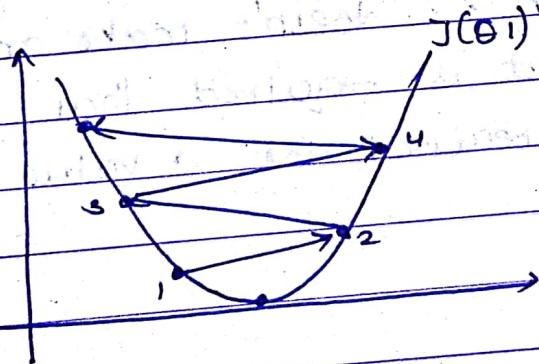
i) α is too small

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



If the α is too small the update by a very small no. so the next point is very near to the previous step as shown in the plot this requires lots of step to reach the global minimum

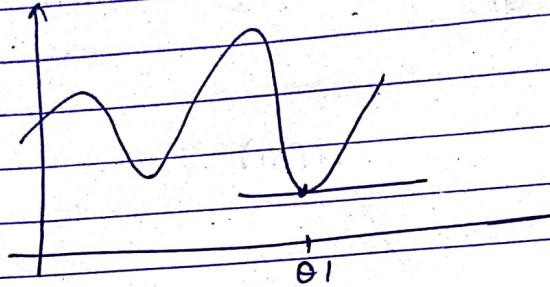
ii) α is too Large



If the α is too large the gradient descent may even overshoot the global minima and fail to converge or even diverge

★ Suppose if the θ_1 is already initialized at the global minima

* Gradient even with



If a slope is drawn to the global minima the slope is 0 ie the derivative is zero as the equation infers.

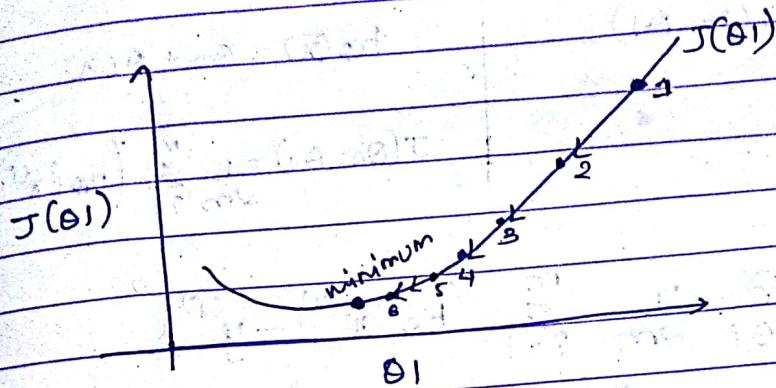
$$\theta_2 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$$\theta_2 := \theta_1$$

This infers that if we already at the global minima any step taken (λ) doesn't make any change. This is actually what is required that once the global minima is reached the θ value should remain unchanged.

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

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



As we approach a local minimum, gradient descent will automatically take smaller steps. So no need to decrease α over time. This infers that if we start θ value from point 1 in the plot & the derivative is larger while α is constant throughout, so the next step point 2 is larger slope at point 2. The derivative is less when multiplied with α it is a smaller value this continues till the local minima is reached. Once the local minima is reached the slope is zero that is derivative is zero and the θ value is fixed. Hence change in α doesn't make any major difference.

* Applying gradient descent algorithm for linear regression Model.

Gradient descent Algorithm

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

Linear regression Model

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

cost function
to be

Batch
↓

each
trainin

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

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

$$\text{Or } j=0: \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\text{Or } j=1: \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

∴ Gradient descent algorithm will be
repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}}$$

$$\theta_1 := \theta_1 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}} x^{(i)}$$

}

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

cost function for linear regression is always going to be a bowl-shaped function (convex function)

Batch Gradient Descent :-

↓
each step of gradient descent uses all of the training examples.

* Linear Regression with multiple variables

multiple features - n

here the features may be

n - no. of features

$x^{(i)}$ - input (feature) of i^{th} training example:

$x_j^{(i)}$ - value of feature j in i^{th} "

ex: size (feet ²)	No. of bedroom	No. of floors	Age of home	Price.
2104	x_1	x_2	x_3	x_4
1416	5	1	45	460
1534	3	2	40	232
852	3	2	30	315
:	2	1	36	178
	:	:	:	:

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 852 \end{bmatrix} \quad x_3^{(2)} = 2$$

previously hypothesis was

$$h_0(x) = \theta_0 + \theta_1 x_1$$

now for multiple variate

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

for the ex

in general

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$x_0 = 1$$

parameter vector

$$\begin{matrix} x = & \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} & \theta = & \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \end{matrix}$$

$$\therefore h_0(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$= \theta^T x.$$

* Gradient descent for Multiple Variants

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

$\theta_0, \theta_1, \dots, \theta_n$ - parameters \Rightarrow $n+1$ dimensional vector

Cost function

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



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

Gradient descent

Repeat i

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

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

OR

learning rate.

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

\Rightarrow gradient descent changes as follows for linear regression for multiple variate

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

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

considering more than 2 features then we have the following simultaneous update for the gradient descent

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

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

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

Feature Scaling

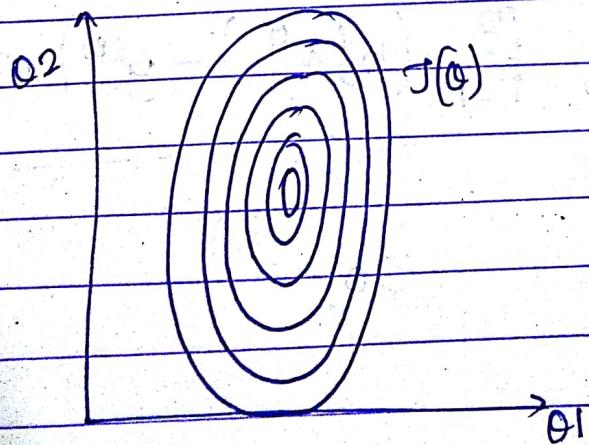
If there are multiple feature and are on similar scale then feature scaling can be used to arrive at a gradient descent.

Consider an example with two feature

$$x_1 = \text{size (0-2000 feet}^2)$$

$$x_2 = \text{number of bedrooms (1-5)}$$

plotting for the above features ignoring θ_0 we have (plotting contours)



As the contours are steep the gradient descent may end up

taking lot of time to reach the global minima.

In these scale the f redefining $x_1 = \dots$

$$a_2 = \dots$$

the contours and when direct path rather than

Feature scaling in the this value is

Normalizing

In these setting the useful thing to do is scale the features that is redefining the x_1 and x_2 as

$$x_1 = \frac{\text{size (feet}^2)}{2000}$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$

then the counters may look much more like circles and when gradient descent is run over this a direct path is found to get the global minima rather than complicating the path

Feature scaling is getting feature approximately in the range of $-1 \leq x_i \leq 1$
this can also be some relatively closer

value like $0 \leq x_1 \leq 3$ ✓

$-2 \leq x_2 \leq 0.5$ ✓

$-0.0001 \leq x_3 \leq 0.0001$ ✗

$-100 \leq x_4 \leq 100$ ✗

Normalization:

To bring data to a same scale

Mean normalization:

Replace x_i with $\frac{x_i - \mu_i}{s_i}$ to make feature to approximately zero mean
(doesn't apply to $\mu_0 = 1$)

$$\text{Eg: } x_1 = \frac{\text{size} - 1000}{2000}$$

$$x_2 = \# \text{bedrooms} - 2$$

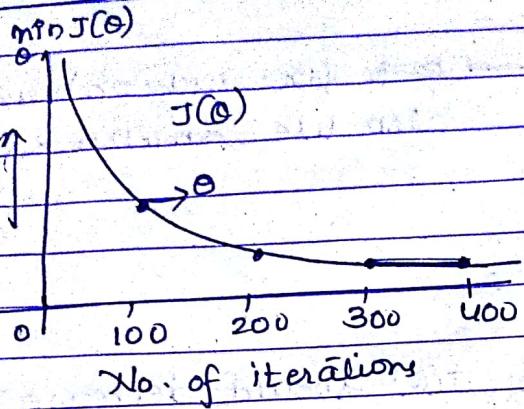
$$-0.5 \leq x_1 \leq 0.5$$

$$-0.5 \leq x_2 \leq 0.5$$

$$x_1 = \frac{x_1 - \mu_1}{s_1} \quad \begin{matrix} \rightarrow \text{avg value of } x \\ \text{training set.} \end{matrix}$$

$$x_2 = \frac{x_2 - \mu_2}{s_2}$$

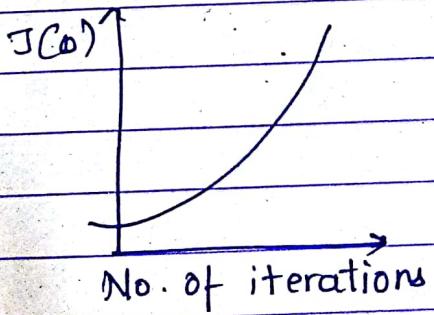
gradient descent II learning rate :



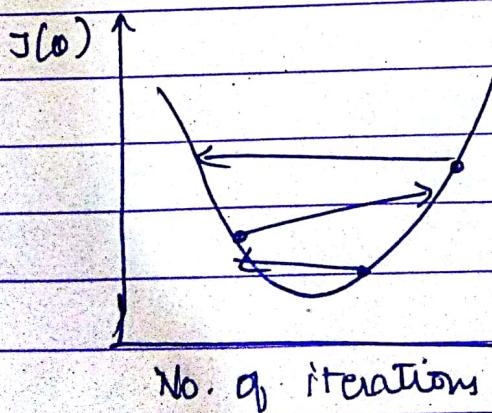
This graph shows the value of θ cost function after certain iterations and $J(\theta)$ should decrease after iterations.

It can be noticed that from 300 - 400 there is no much decrease in the gradient this implies that there is no ~~shar~~ major change in the θ and ~~this~~ at ~~many~~ these iteration there may be convergence of the gradient descent.

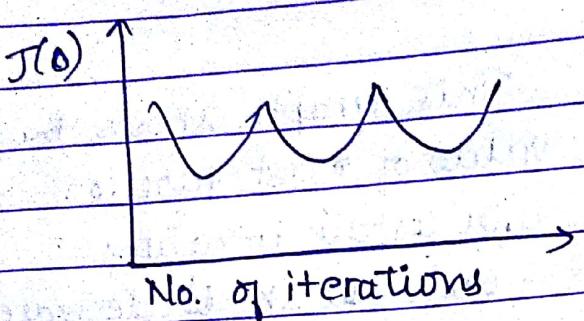
* Making sure the gradient descent is working properly



here $J(\theta)$ is increasing this notifies that the gradient descent is not working
 $\therefore \alpha$ should be smaller



- if α is too large the gradient descent keeps on over shooting.
 \therefore use ~~a~~ is smaller α in such cases.



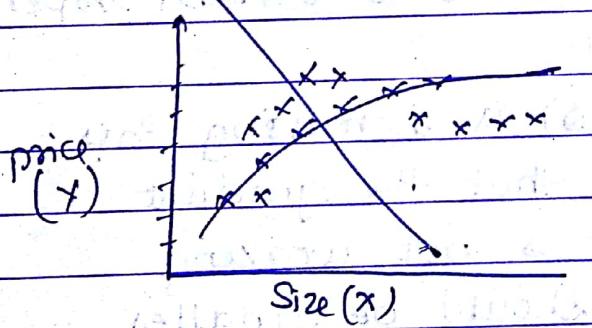
→ to fix such problem
also use smaller α

inference from the above

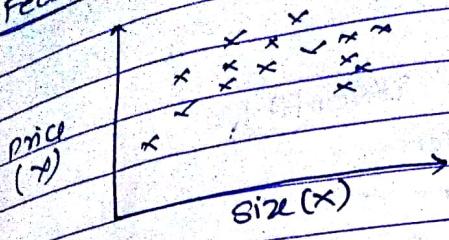
- For sufficiently small α , $J(\theta)$ should decrease on every iteration
- But if α is too small, gradient descent can be slow to converge

* Features & Polynomial Regression

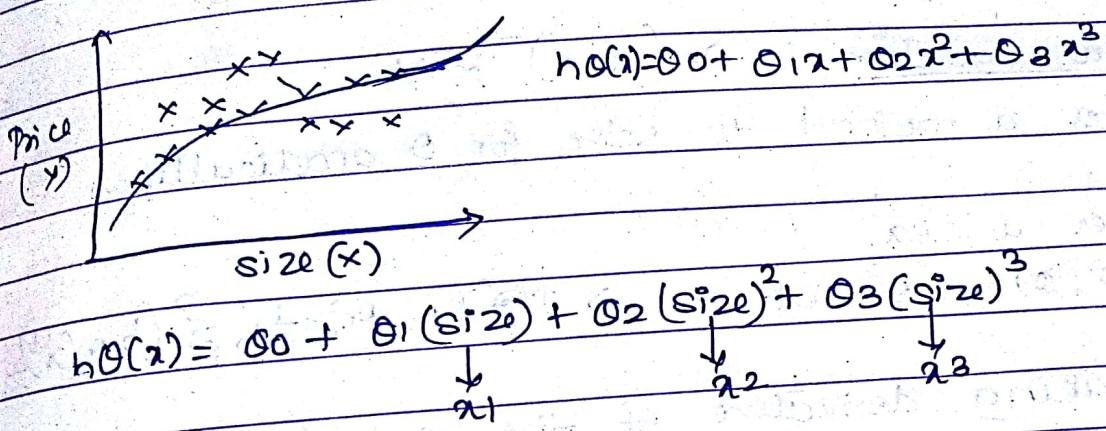
Consider a ex



Features and polynomial regression

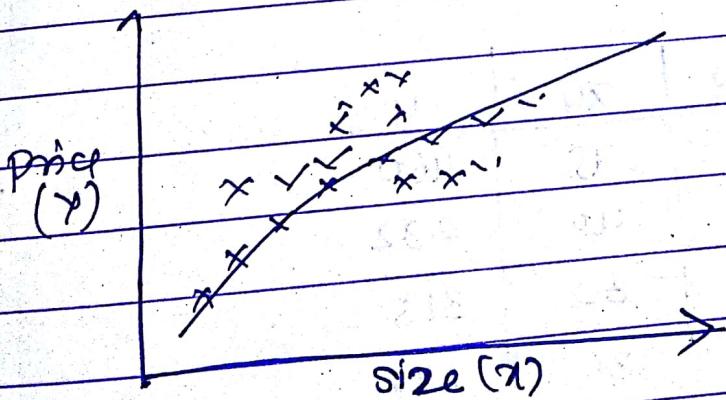


A straight line cannot fit into the data
and a quadratic function such as $\theta_0 + \theta_1 x + \theta_2 x^2$
but eventually goes down
therefore choosing 3rd order function



Another choice to fit the data is

$$h_{\Theta}(x) = \theta_0 + \theta_1 \text{size} + \theta_2 \sqrt{\text{size}}$$



~~Normal Equations~~: Non invertibility

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

$X^T X$ is sometimes non invertible
to avoid or overcome then

→ redundant features

→ If there many features ($m < n$)
delete some features or regularize

Normal Equations: $\theta = (X^T X)^{-1} X^T Y$

gives a method to solve for θ analytically

consider a ex.

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

taking derivative of the cost function

and solving for the derivative we arrive
at θ values that take to global minima

Ex:-

$m=4$

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



adding $n+1^{\text{th}}$ training Sample.

forming matrix for x and y and solving for
 $\theta = (x^T x)^{-1} x^T y$ give the θ value that ends up
reaching the global minima.

$$x = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad y = \begin{bmatrix} 460 \\ 232 \\ 15 \\ 178 \end{bmatrix}$$

$$x = m \times (n+1)$$

\downarrow
n. dimensions

Generalizing the equations

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

n features

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

$$(x) \text{ design matrix} = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(m)})^T - \end{bmatrix}$$

$$x = m \times (n+1)$$

LP-5

Gradient Descent

→ disadvantage: need to choose α
Needs many iteration

→ advantage:

: work well even
when n is large

Normal Equations

adv: No need to choose α
Don't need to iterate

disadvantage:

* Need to compute
 $(X^T X)^{-1}$
* slow if n is large

Logistic regression

earlier classification (linear classification) the prediction y may have two values i.e. $y \in \{0, 1\}$

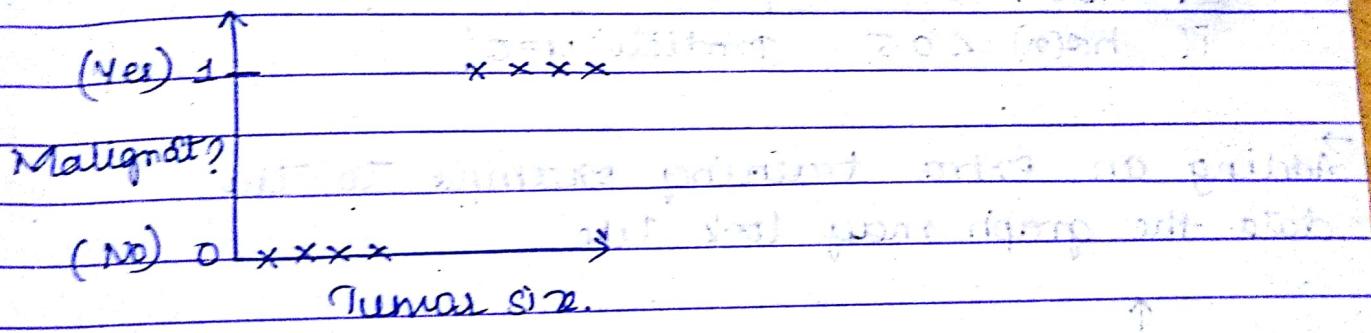
0 - Negative class (ex: benign tumor)

1 - +ve class (ex: malignant tumor)

(logistic regression has binary classification and multivariate classification as well)

drawback of linear regression

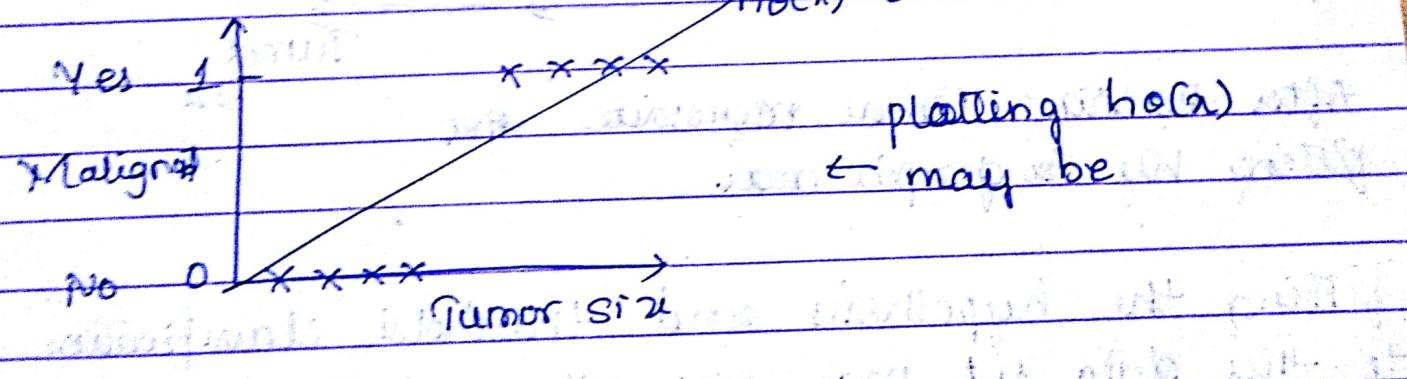
Ex:



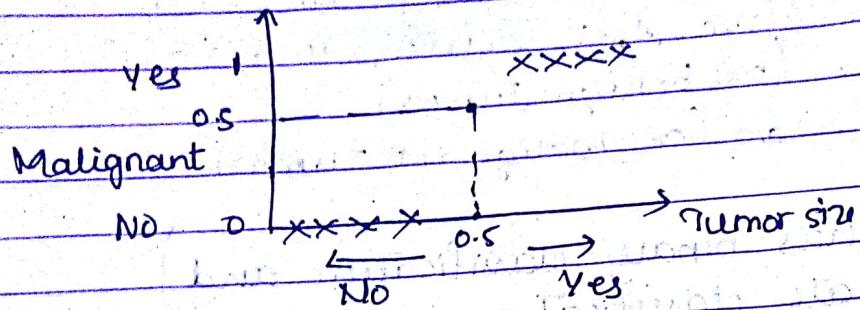
Considering a hypothesis function (used earlier)

$$h_0(x) = \theta^T x$$

(linear function)



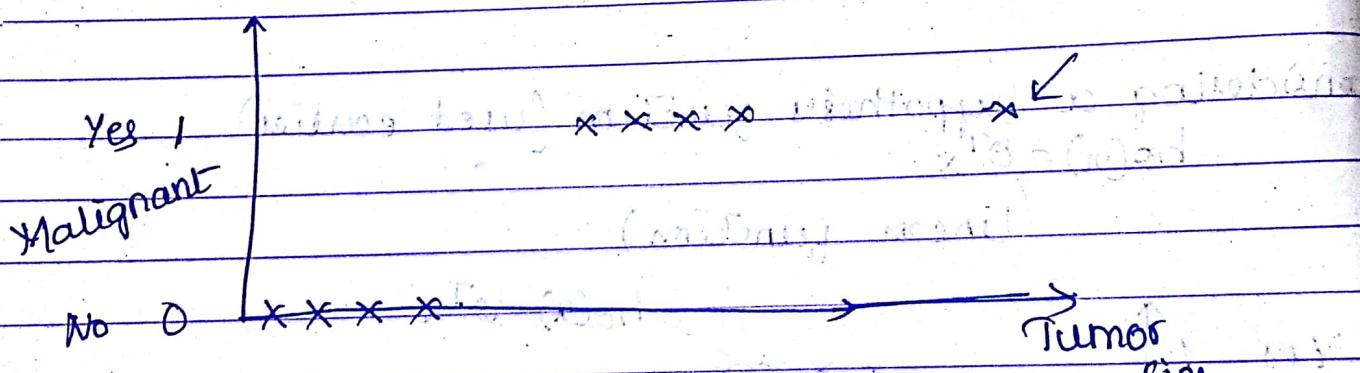
Considering threshold.



Threshold classification output $h_0(x) @ 0.5$

if $h_0(x) \geq 0.5$ predict 'y = 1'
if $h_0(x) < 0.5$ predict 'y = 0'

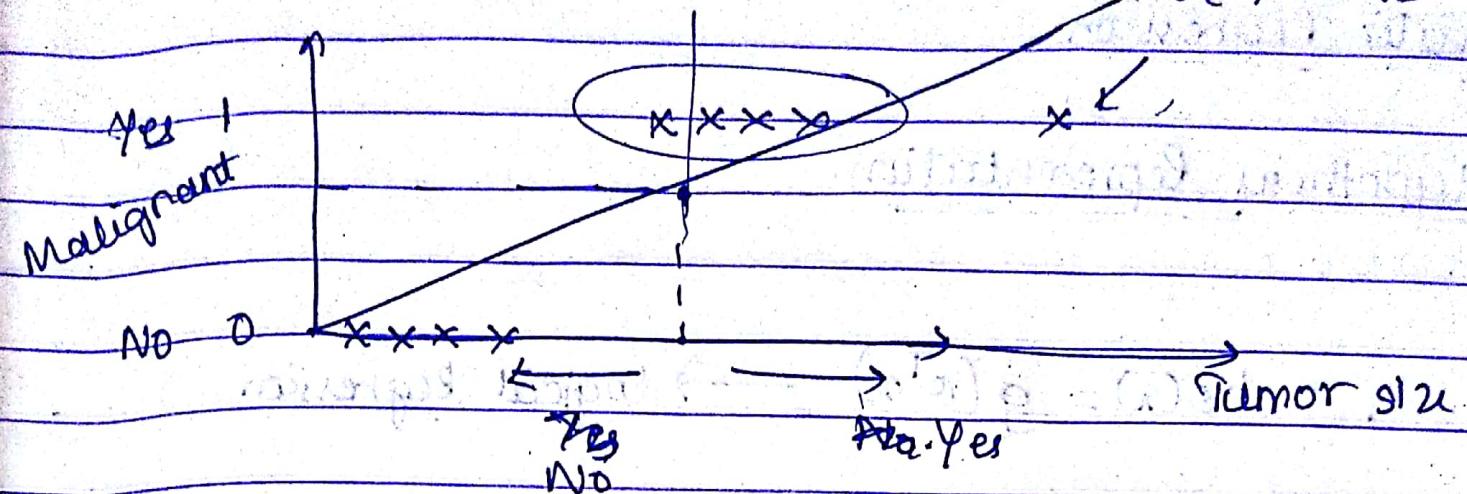
→ Adding an extra training example to the data the graph may look like



After running linear regression the fitting line for graph was

fitting the hypothesis and threshold classification to this data set may vary the whole prediction

$$h_0(x) = \phi^T x.$$



The training sets on Yes(1) line are all the ex- with Malignant tumor but when the linear regression is run on this data set the best fitting of hypothesis changes and the threshold classification also changes as it can be seen from the above diagram that few of the malignant tumor data ex who lie on the either side of begin tumor. ∴ linear regression is not preferred for classification problem.

* The y takes values 0 and 1 but $h_0(x)$ can be >1 or <0

this is constraint for fitting the curve.

where as if $\phi^T x$ $h_0(x)$ is 0 or 1 is logistic Regression

logistic regression

Hypothesis Representation

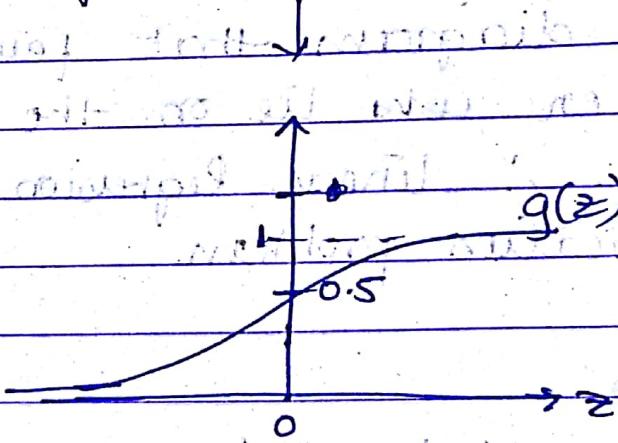
e.g. from $h_{\theta}(x) = g(\theta^T x) \rightarrow \text{logistic regression}$

for logistic

$$h_{\theta}(x) = g(\theta^T x)$$

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

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



Interpretation of Hypothesis Output.

$h_0(x)$ = estimated probability that $y=1$ on input x

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

then $h_0(x) = 0.7$.

Suppose the patient have some tumor size and their feature hypothesis is 0.7. This tells that for the patient with features x the probability that the patient is having tumor is 70% (0.7).

Mathematically writing.

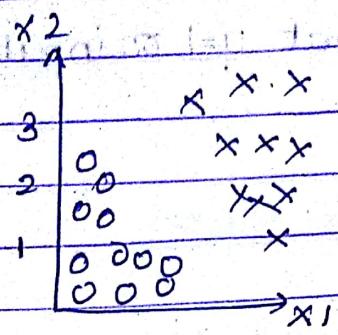
$$h_0(x) = P(y=1|x; \theta) \quad \text{"probability that } y=1 \text{, given } x, \text{ parameterized by } \theta"$$

y can take either 0 or 1.

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

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

Decision Boundary



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

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \theta_0 = -3 \quad \theta_1 = 1 \quad \theta_2 = 1$$

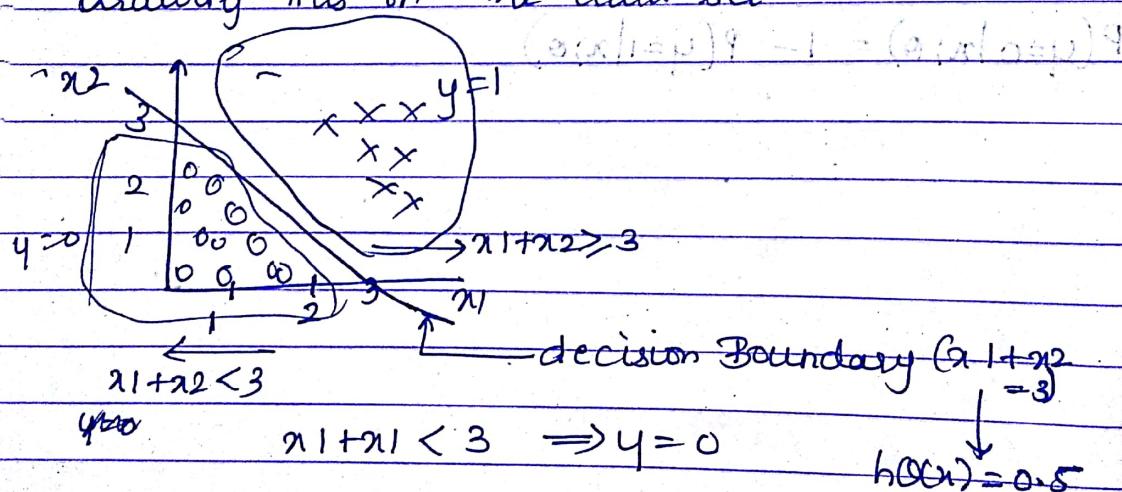
Predict $y = 1$ if $x_1 + x_2 \geq 0$ ($(x_1, x_2) = (0, 0)$)

$$x_1 + x_2 \geq 0$$

considering $x_1 + x_2 = 3$

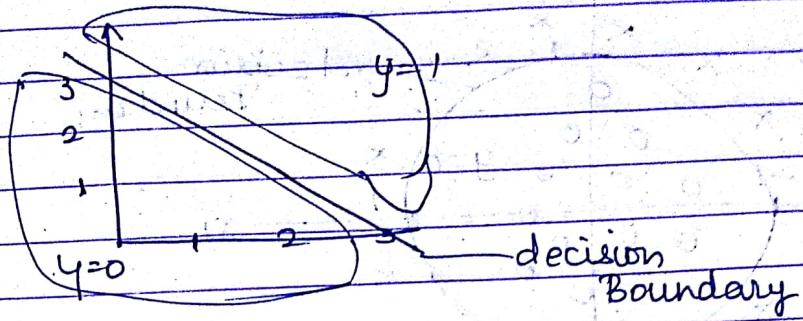
this is a straight line

drawing this on the data set

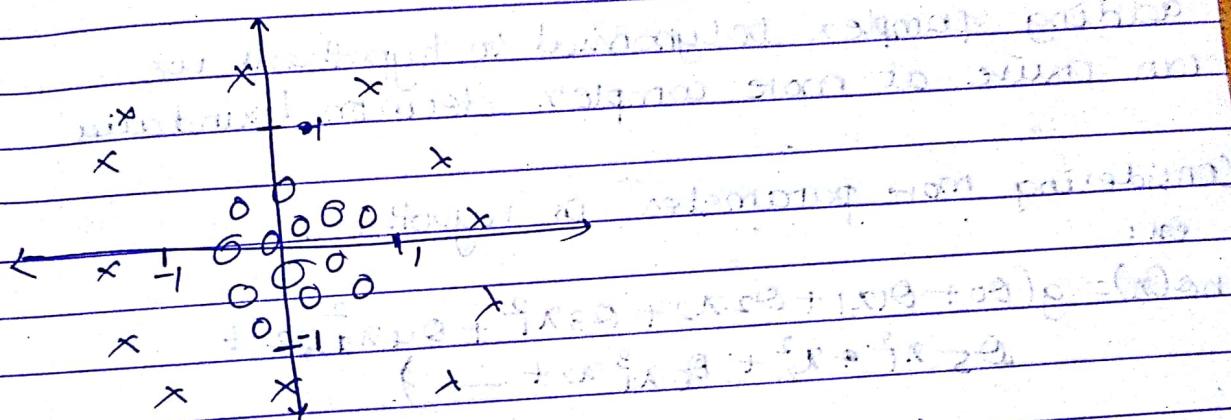


Decision boundary is the line where that separates where the hypothesis predicts y is exactly equal to 0 and exactly equal to 1

Decision boundary is the property of hypothesis and the parameter and not the property of the data set



Non linear decision Boundaries



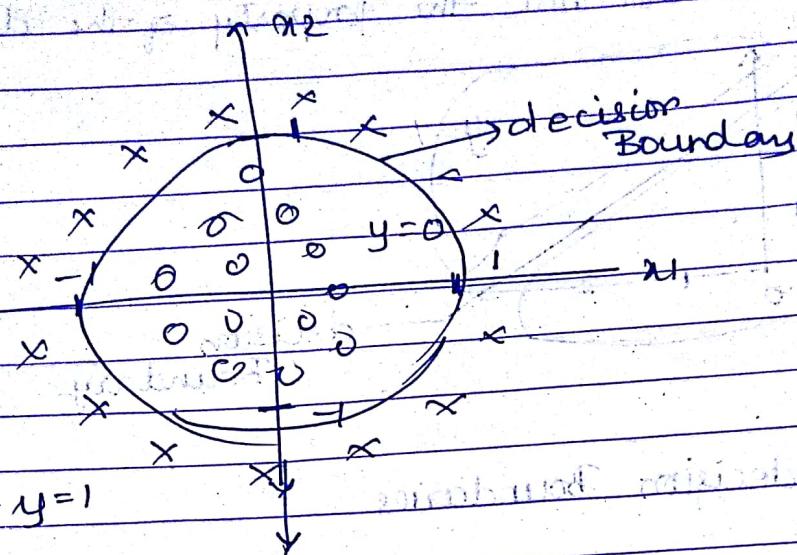
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \begin{aligned} \theta_0 &= -1 \\ \theta_1 &= 0 \\ \theta_2 &= 0 \\ \theta_3 &= 1 \\ \theta_4 &= 1 \end{aligned}$$

Predict $y=1$ if $-1 + x_1^2 + x_2^2 \geq 0$

$$x_1^2 + x_2^2 \geq 1$$

$$x_1^2 + x_2^2 = 1 \rightarrow \text{equation of circle}$$



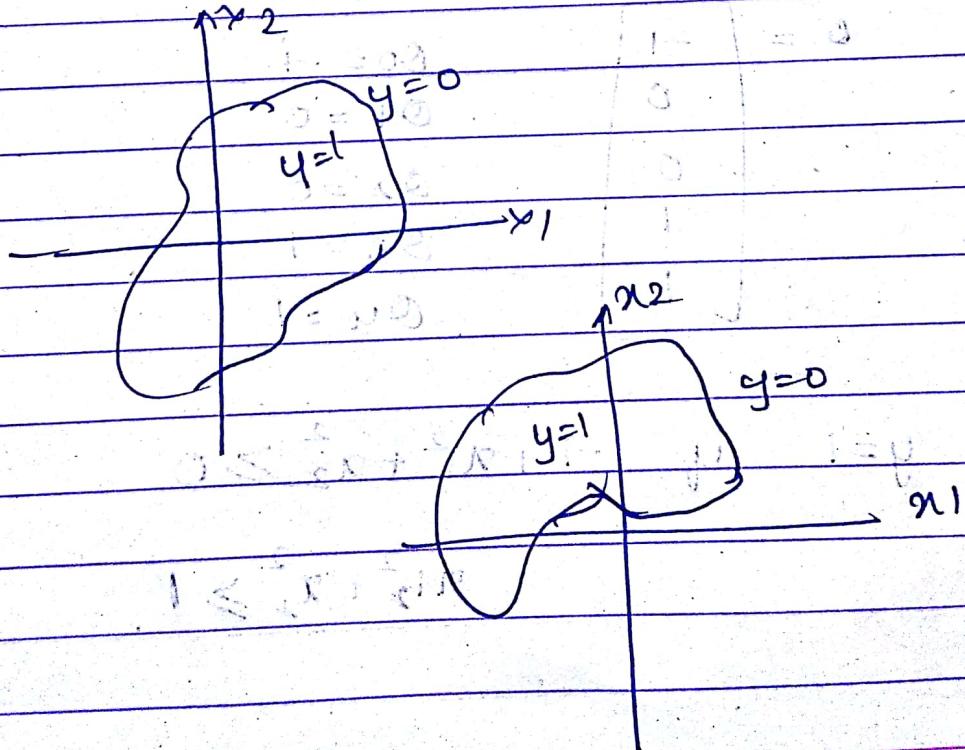
adding complex polynomial to hypothesis we can arrive at more complex decision boundaries

Considering more parameters in hypothesis

e.g:

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

The decision boundary may end up looking like (x)



Fitting Data

tn:

Training set

In example

$h_{\theta}(x)$

choose

Fitting Data to Logistic Regression

in:

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

m examples

$x \in$	$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$
---------	---

$$x_0 = y, y \in \{0, 1\}$$

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

choosing parameter θ to (cost function)

Cost function

Linear Regression

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

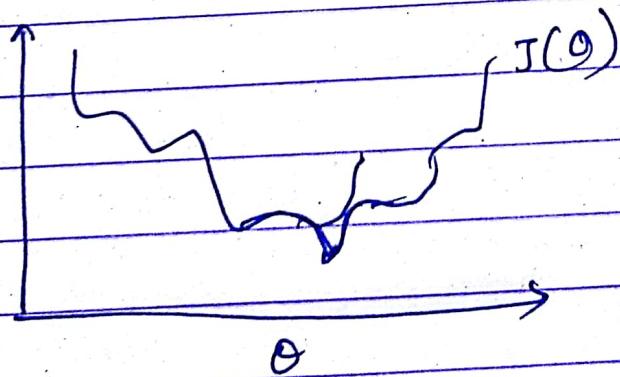
$$\text{cost}(h_\theta(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

Considering same cost function for Logistic Regression

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

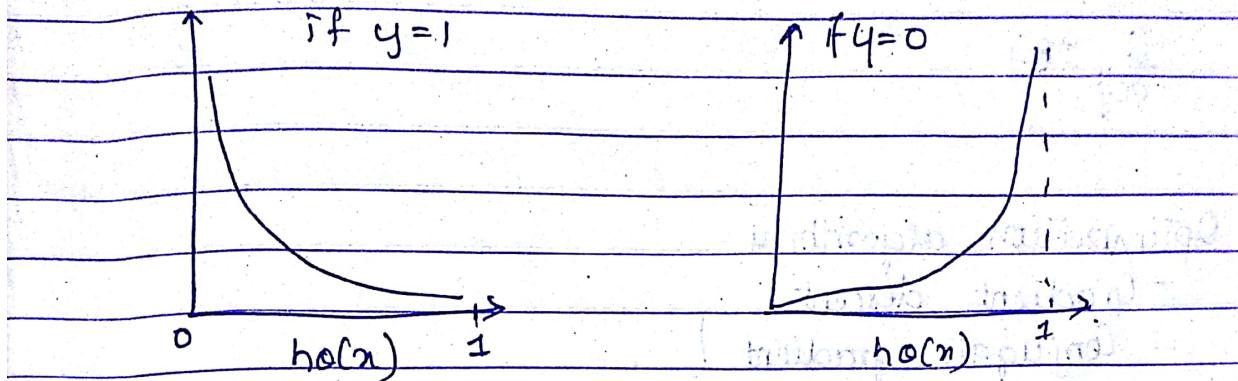
(non linear function)

The resulting gradient may look like a non-convex curve



Cost function for logistic regression

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



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

can also be written as

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

To fit parameters:

$$\min_{\theta} J(\theta) \rightarrow \text{get } \theta$$

repeat q

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

(simultaneously update all θ_j)

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

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

* Optimization Algorithms

given θ we need that can compute

$$J(\theta)$$

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

Optimization algorithms

- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages

- No need to manually pick α
- Often faster than gradient descent

disadvantages

- More complex.

Mult class Classification:

Ex:

Email foldering / tagging: Work, Friends, Family + Hobby

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$y=1$$

$$y=2$$

$$y=3$$

$$y=4$$

Medical diagnosis: Not ill, cold, flu

$$\downarrow \quad \downarrow \quad \downarrow$$

$$y=1$$

$$y=2$$

$$y=3$$

Weather: sunny, cloudy, rain, snow

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

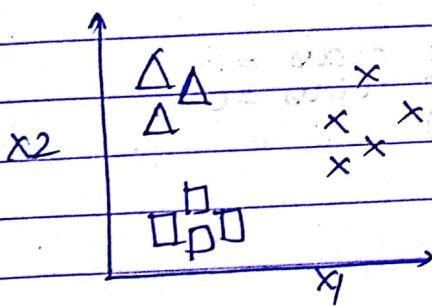
$$y=1$$

$$y=2$$

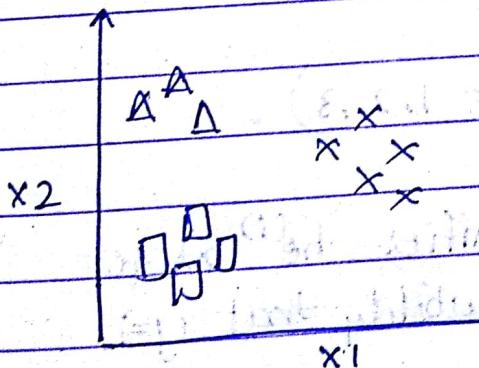
$$y=3$$

$$y=4$$

Ex for data set



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

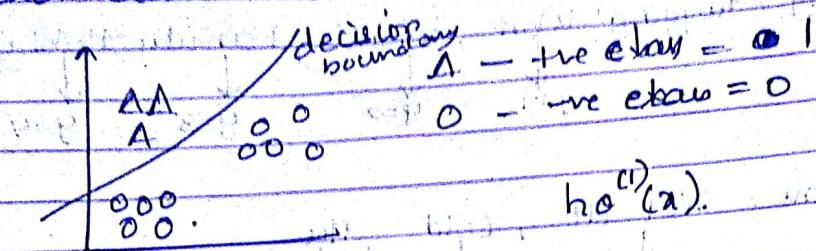


Class 1 = □

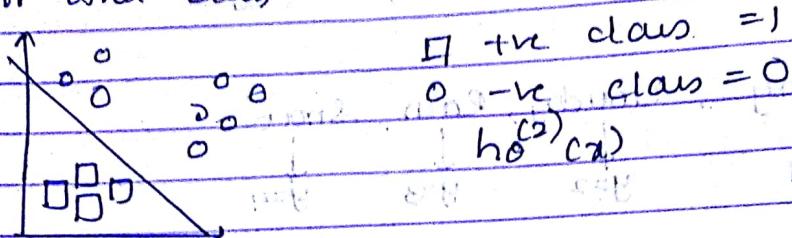
2 = □

3 = x

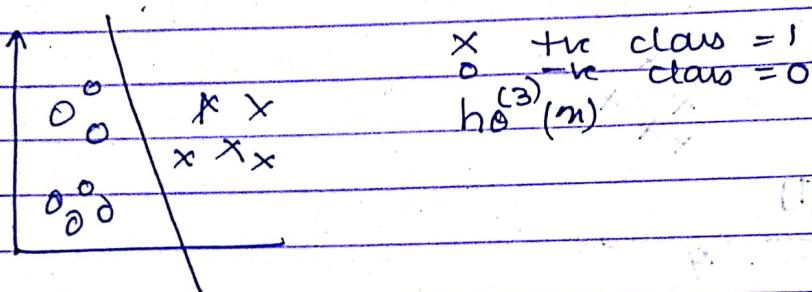
Considering class we create new training set
where class 2 & 3 are assigned to -ve class
and class 1 is +ve class



try for other class



try for 3rd class



To summarize we fit 3 classifier then
find the probability that $y=1$ give for any
 x values

$$h_{\theta}^{(i)}(x) = P(y=i|x; \theta) \quad (i=1, 2, 3)$$

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

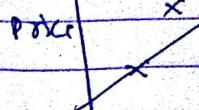
On a new input x to make a prediction, pick the
class i that maximizes

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

The Problem

Linear Regre

Ex: Housing



Dot+

Under

price

Dot+

ju

Overselli

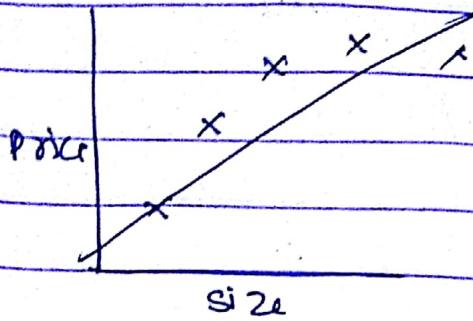
hypothesis

but

The Problem of overfitting

Linear Regression

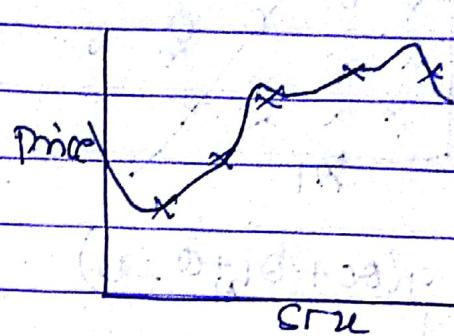
Ex: Housing Prices



$$\theta_0 + \theta_1 x$$

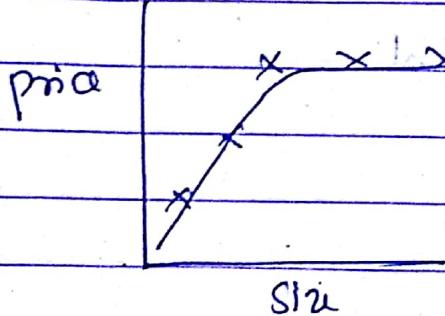
Underfit

High bias



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

Overfit - high variance

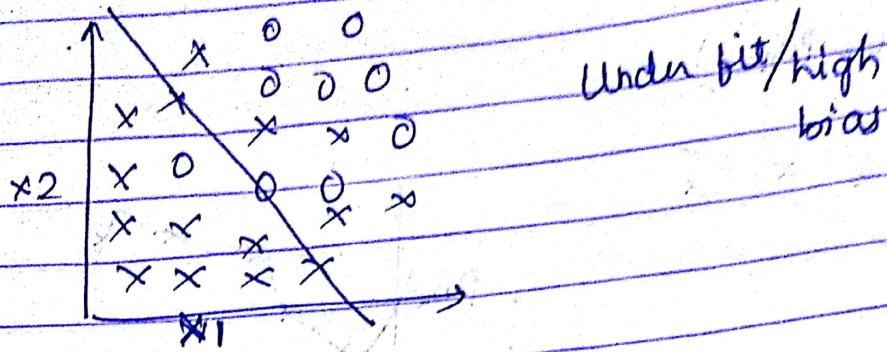


$$\theta_0 + \theta_1 x + \theta_2 x^2$$

just right

Overfitting: if we have too many features, the learned hypothesis may fit the training set very well but fail to generalize to new examples

logistic regression



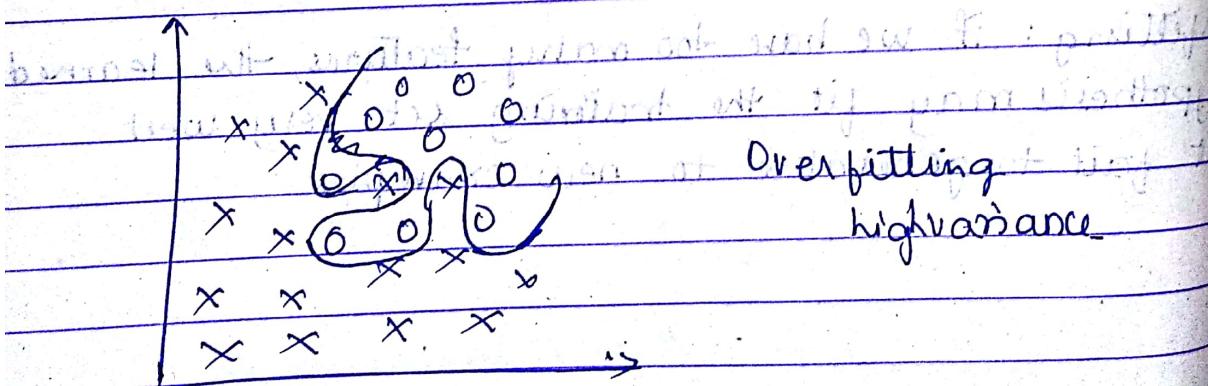
$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

sigmoid function

nonlinear output - sigmoid

just right

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

Addressing Overfitting:-

disadvantage:
ignoring of important
few features

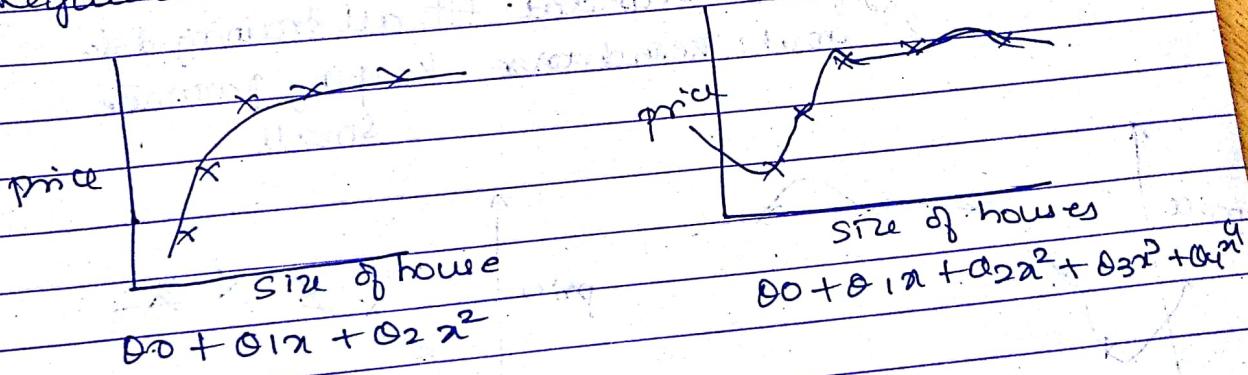
1. Reduce Number of feature

- Manually select which features to keep
- Model selection algorithm

2. Regularization

- Keep all the features but reduce magnitude/ values of parameters θ_j
- Works well when we have a lot of features each of which contributes a bit to predicting y

Regularization



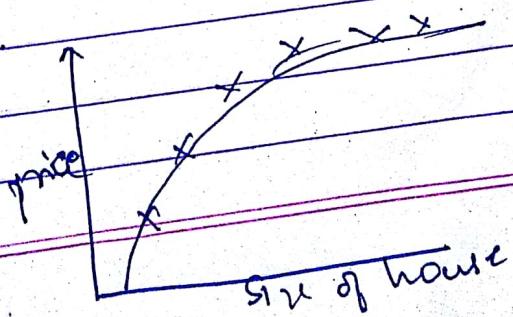
Suppose θ_3 & θ_4 are made very small.

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

some Random huge value!

$\therefore \theta_3 \approx 0$ and $\theta_4 \approx 0$.

\therefore we end up having quadratic eqⁿ with negligible value of θ_3 & θ_4 .



small value for parameter

- "simple" hypothesis

- less prone to overfitting

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

→ Regularized cost function

regularization term

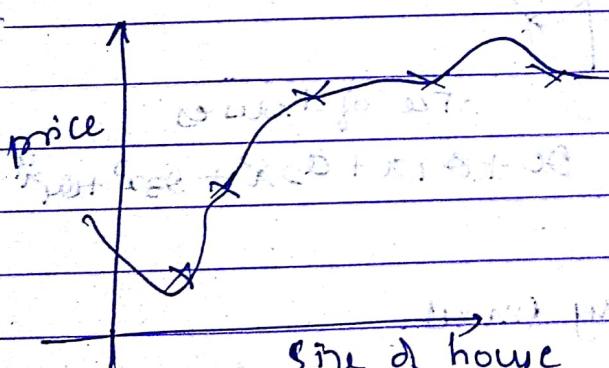
regularization parameter

controls trade off betⁿ

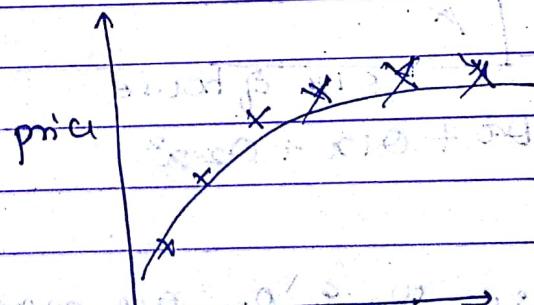
two goals

1st goal: trained to fit all training data

2nd goal: regularize keeping parameter small.



without regularization



with regularization

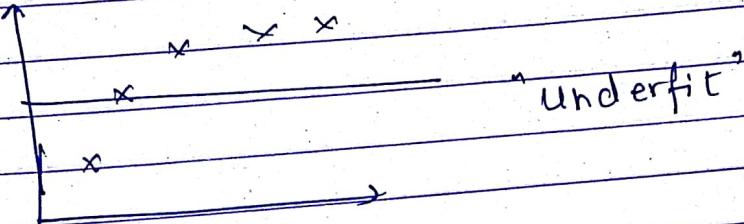
if λ is set too large, then

$h_{\theta}(x)$

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

become $\equiv 0$

$$\therefore h_{\theta}(x) = \theta_0.$$



regularized linear regression
gradient descent

Repeat {

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

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

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

$$\frac{1 - \alpha \lambda}{m} \leq 1$$