

Machine Learning

Arthur Samuel:-

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

Tom Mitchell:-

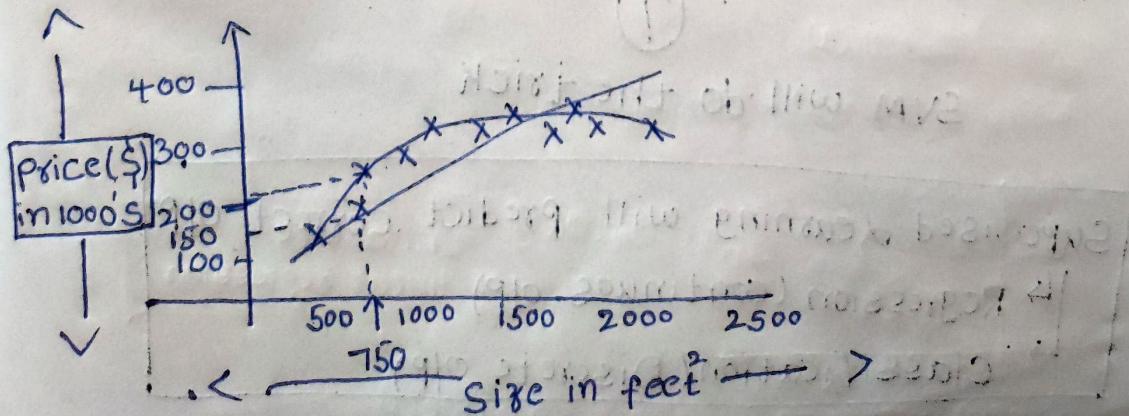
A. Computer program is said to learn from experience E with respect to some task T & some performance measure P, if its performance on T, as measured by P, improves with experience

E.

Machine Learning Algorithms

- Supervised learning
- Unsupervised learning
- Reinforcement learning, Recommender Systems

Supervised learning

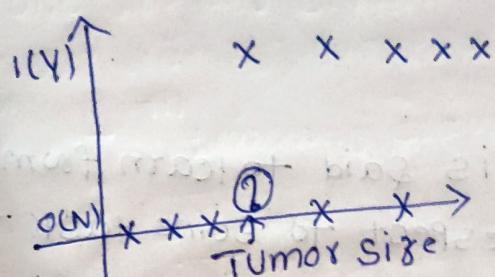


• Regression problem (Predict continuous valued O/P)

Here in this case its (price)

④ Breast Cancer (malignant, benign) [classification problem]

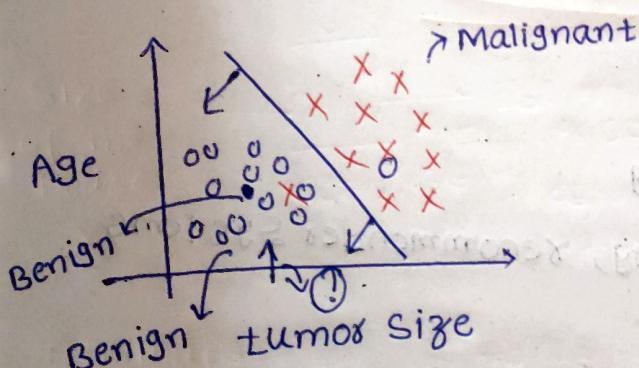
e.g:-



Classification
Discrete valued
O/P (0 or 1)

Here we considered "tumor size" as an attribute
so this is just for one attribute:

Considering more than 1 attribute



⇒ If there are infinite no of features
(?)

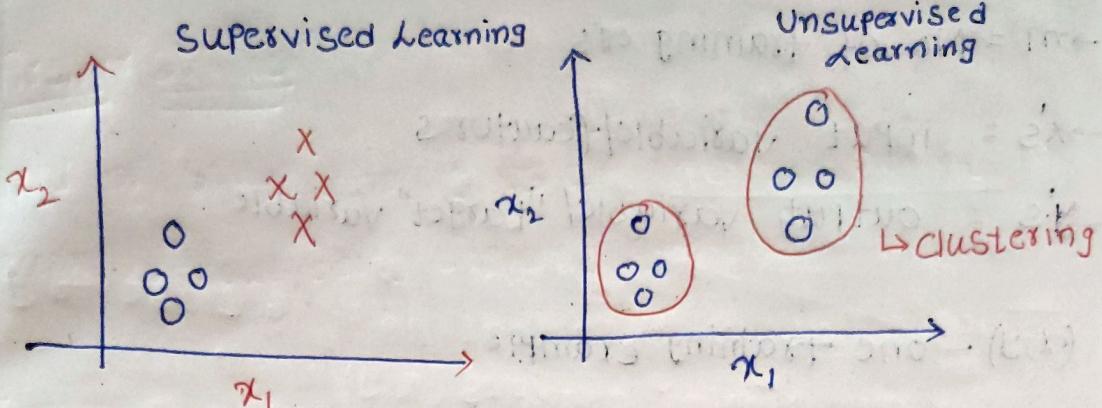
SVM will do the trick

Supervised learning will predict correct O/P

↳ Regression (continuous O/P)

↳ Classification (discrete O/P)

Unsupervised Learning

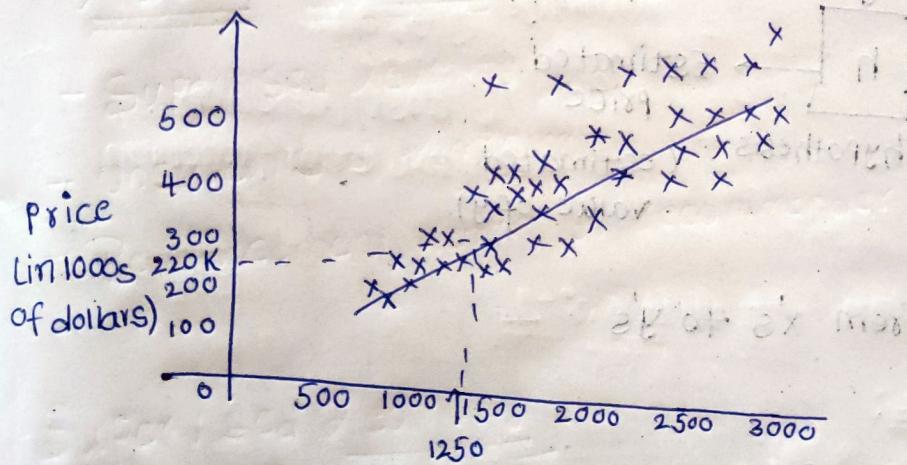


[Labels are present]

[No labels]

Given a dataset & you should find the structure b/w the data points

Model Representation



Supervised learning

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

Regression (Problem)

Predict real-valued o/p
↳ continuous

Notation

$\rightarrow m$ = No. of training eq's

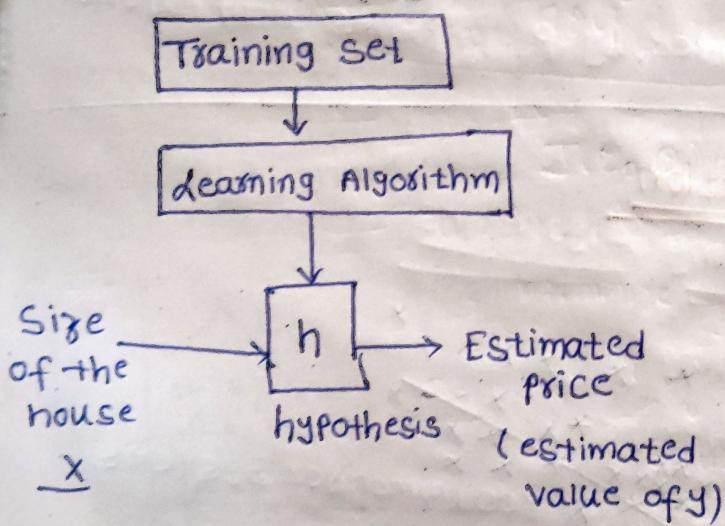
$\rightarrow x$'s = input variable/features

$\rightarrow y$'s = output variable/"target" variable

(x, y) - one training example

$(x^{(i)}, y^{(i)}) \rightarrow i^{\text{th}}$ training example

How we build a model

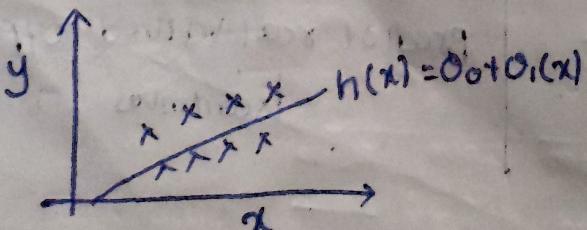


h maps from x 's to y 's

Representation of h ?

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

shorthand: $h(x)$



univariate
linear regression

linear regression with one variable(x)

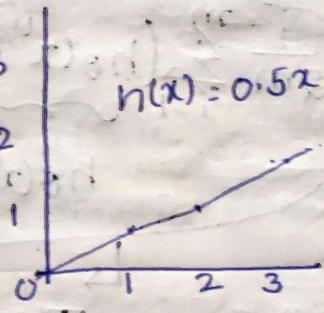
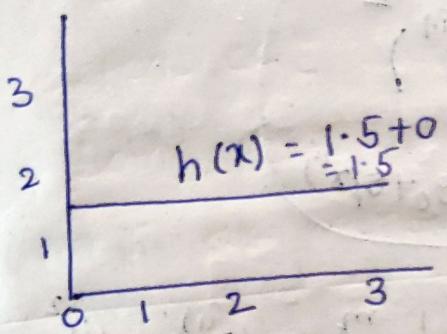
Cost function [This will decide the best fit line in our data]

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

θ_i 's : Parameters

How to choose θ_i 's?

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



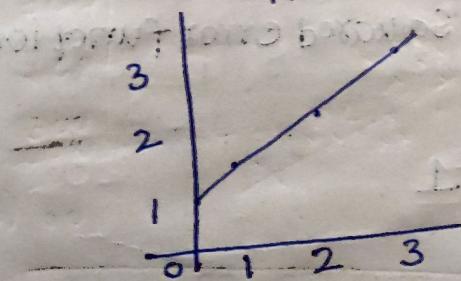
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$

$$\theta_0 = 0$$

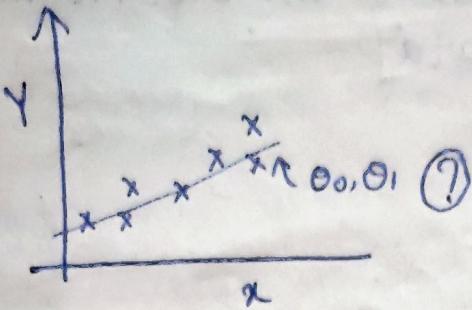
$$\theta_1 = 0.5$$

$$h(x) = 1 + 0.5x$$



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$



Idea: Choose θ_0, θ_1 so that $h_0(x)$ is close to y for our training examples (x, y)

minimize
 θ_0, θ_1

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

$$h_0(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

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

minimize $J(\theta_0, \theta_1)$

θ_0, θ_1 ✓ cost function

squared error function

Cost function - Intuition - 1

Hypothesis:

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

Parameters

$$\theta_0, \theta_1$$

$$\cancel{\underline{h(x)}}$$

Cost function

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

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

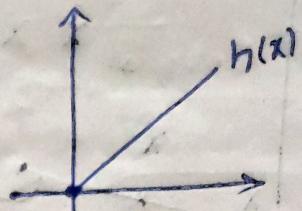
Simplified form

$$h_{\theta}(x) = \theta_1 x \quad (\theta_0 = 0)$$

$$\theta_0 = 0$$

θ_1

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



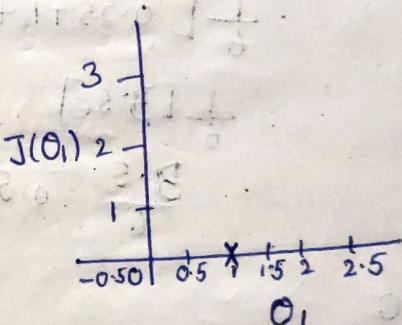
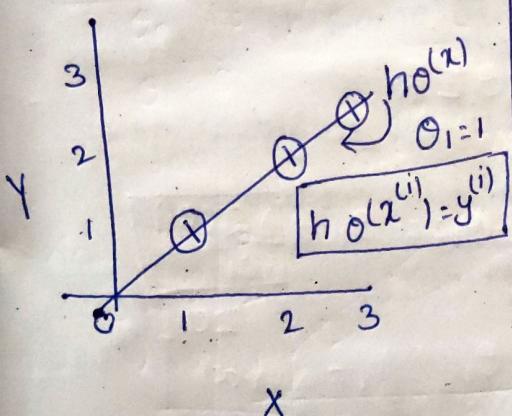
minimize $J(\theta_1)$

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

(for fixed θ_1 , function
of x)

$$J(\theta_1)$$

(function of Parameter θ_1)



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

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

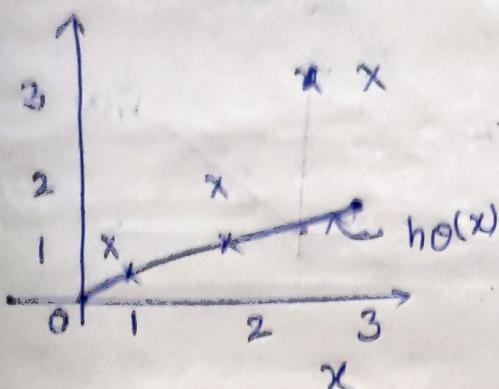
$$= \frac{1}{2m} [0^2 + 0^2 + 0^2]$$

$$J(\theta_1) = 0$$

$$J(1) = 0$$

$$\theta_1 = 0.5$$

$$h_0(x) = 0.5x$$



$$J(\theta)$$

3.5
3
2.5
2
1.5
1
0.5
0

$$J(0)$$

$$J(0.5) = \frac{1}{8m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$

$$= \frac{1}{8m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

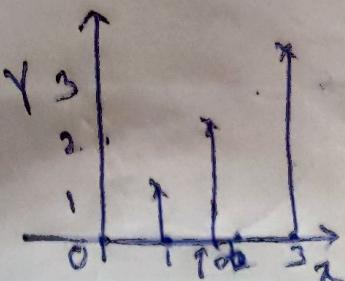
$$= \frac{1}{6} [0.25 + 1 + 2.25]$$

$$= \frac{1}{6} [3.50]$$

$$= \frac{3.5}{6} = 0.58$$

$$\theta_1 = 0$$

$$h_0(0) = 0(2) = 0$$

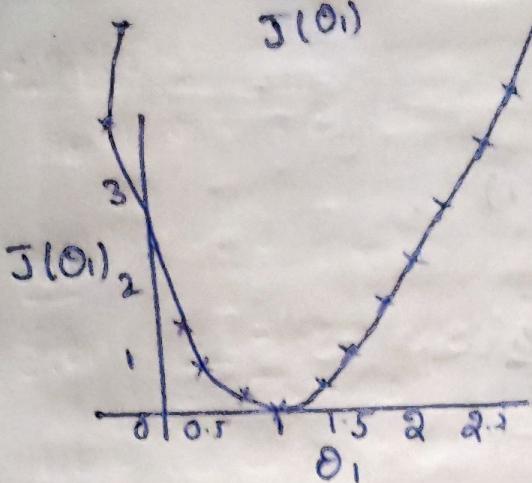


2.5
2
1.5
1
0.5
0

$$J(\theta_1)$$

$$\therefore J(0) = \frac{1}{6} \sum_{i=1}^m [h_0(2^{(i)}) - y^{(i)}]^2$$

$$= \frac{1}{6} [1+4+9] = \frac{14}{6} = 2.33$$



Gradient Descent : [Minimize the cost function & other facts as well]

$$\Rightarrow J(\theta_0, \theta_1)$$

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

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

Gradient descent algorithm

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

Simultaneously update

θ_0 & θ_1

Correct

Correct Simultaneous update

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

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

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

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

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

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

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

$$\theta_0 = 1 \text{ & } \theta_1 = 2$$

$$\theta_j := \theta_j + \sqrt{\theta_0 \theta_1} \quad (\text{for } j=0 \text{ & } j=1)$$

$$\theta_0 := \theta_0 + \sqrt{2}$$

$$\boxed{\theta_0 := 1 + \sqrt{2}}$$

$$\theta_1 := \theta_1 + \sqrt{\theta_0 \theta_1}$$

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

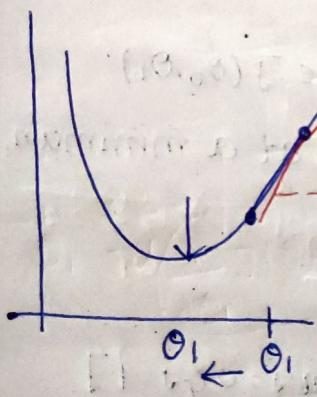
Gradient descent algorithm

Repeat until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \quad [\text{simultaneously update } \theta_0 \text{ & } \theta_1]$$

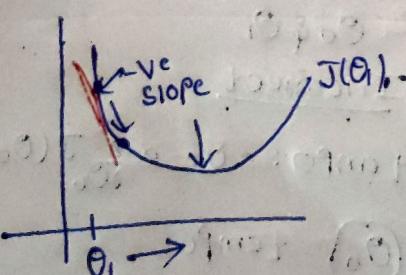
Learning rate derivative

eg:-



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

$$\theta_1 := \theta_1 - \alpha (+ve)$$

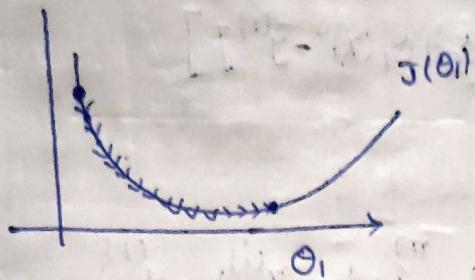


$$\theta_1 := \theta_1 - \alpha (-ve)$$

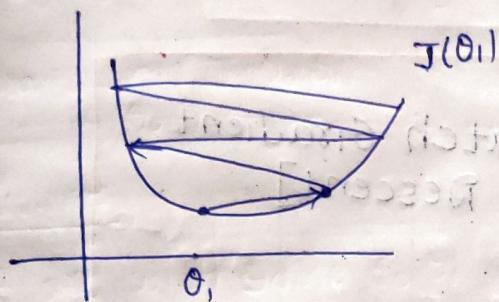
$$\theta_1 := \theta_1 - \alpha (-ve)$$

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

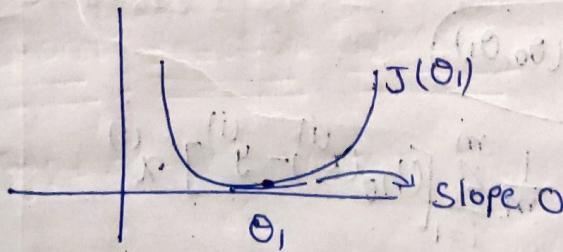
α is too small, gradient descent can be slow



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



If θ_1 is at local minima then what happens?



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

$$\therefore \theta_1 := \theta_1 - \alpha(0)$$

$$\theta_1 := \theta_1 \quad \boxed{\text{Remains Unchanged}}$$

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

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

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

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

Gradient algorithm [Batch Gradient Descent]

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

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

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

$$\xrightarrow{\text{Partial derivative}} \frac{1}{m} \sum_{i=1}^m \left[(h_\theta(x^{(i)}) - y^{(i)}) \right] x^{(i)}$$

Linear Algebra

Matrix Addition:

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

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix}_{3 \times 2} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \end{bmatrix}_{2 \times 2} = [?] \quad \text{X} \quad \text{Can't perform; Should be of square matrices}$$

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

$$\begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} / 4 = \begin{bmatrix} 1 & 0 \\ 3/2 & 3/4 \end{bmatrix} \quad \checkmark$$

Combination of operands

$$3 \times \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} - \begin{bmatrix} 3 \\ 0 \\ 2 \end{bmatrix} / 3 = \begin{bmatrix} 3 \\ 12 \\ 6 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ \frac{2}{3} \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ 12 \\ \frac{31}{3} \end{bmatrix} \quad \checkmark$$

Matrix multiplication

e.g:-

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

$\overset{3 \times 2}{\text{}} \quad \overset{2 \times 1}{\text{}}$

$\hookrightarrow \text{Res: } 3 \times 1$

$$\begin{bmatrix} 16 \\ 4 \\ 7 \end{bmatrix} \quad \overset{3 \times 1}{\text{}}$$

House sizes

2104

1416

1534

852

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

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \begin{bmatrix} -40 \\ 0.25 \end{bmatrix} = \begin{bmatrix} 0 \\ -40 + 2104 \times 0.25 \\ -40 + 1416 \times 0.25 \\ -40 + 1534 \times 0.25 \\ -40 + 852 \times 0.25 \end{bmatrix}$$

$\overset{h_0(2104)}{\text{}} \quad \overset{h_0(1416)}{\text{}}$

2104, 1416, 1534, 852 \Rightarrow no. of houses

Multivariate Linear Regression

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

For convenience of notation define $x_0 = 1$

$$x_0^{(i)} = 1$$

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

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

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

$$\boxed{h_{\theta}(x) = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \cdot \begin{bmatrix} x_0 & x_1 & x_2 & \dots & x_n \end{bmatrix} = \theta^T x}$$

$\Theta \in \mathbb{R}^{n+1}$
 $\begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix}$

$\Theta \in \mathbb{R}^{n+1}$

$$\boxed{h_{\theta}(x) = \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \\ (n+1)x_1 \end{bmatrix} = \begin{bmatrix} \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n \end{bmatrix}}$$

Gradient descent for multiple features

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

Parameters: $\theta_0, \theta_1, \dots, \theta_n$

Assume a vector θ

Cost function:-

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

Gradient descent

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

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j}$$

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

Gradient Descent in Practice 1 - Feature Scaling

Scaling

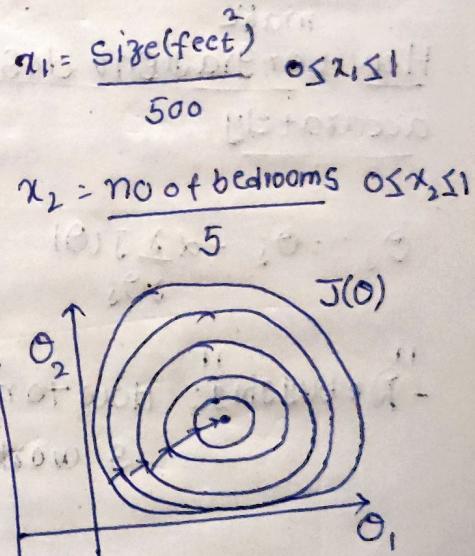
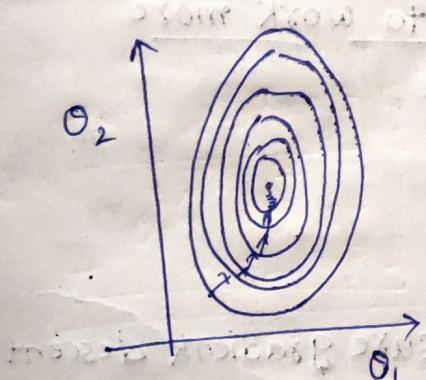
Feature Scaling

Idea: Make sure features are on a similar scale

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

$x_2 = \text{no of bedroom } (1\text{-}5)$

If you implement with these then gradient descent will take huge time to reach global minima



Feature Scaling

Get every feature approximately a $-1 \leq x_i \leq 1$

$$x_0 = 1$$

$$0 \leq x_1 \leq 3 \checkmark$$

$$-2 \leq x_2 \leq 0.5 \checkmark$$

$$-100 \leq x_3 \leq 100 X$$

$$-0.001 \leq x_4 \leq 0.001 X$$

$$-3 + 0.3 \checkmark$$

$$-\frac{1}{3} + \frac{1}{3} \checkmark$$

Mean Normalization

[Do not apply to $x_0 = 1$]

$$x_1 = \frac{\text{Size} - \mu_1}{2000}$$

$$x_2 = \frac{\text{bedrooms} - \mu_2}{5}$$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{\sigma_1}; \mu \text{ avg value of } x \text{ in training set}$$

$$x_2 \leftarrow \frac{x_2 - \mu_2}{\sigma_2}; \sigma \text{ Range}$$

approx: $-0.5 \leq x_1 \leq 0.5$ & $-0.5 \leq x_2 \leq 0.5$

[Should need not be exact]

make
How to gradient descent to work more accurately

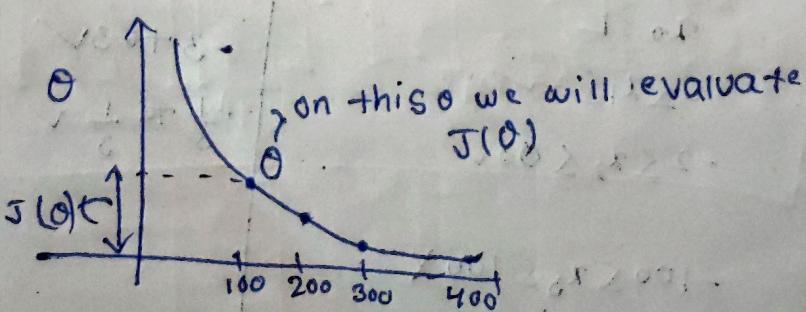
$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

"Debugging": How to make sure gradient descent is working correctly

- How to choose learning rate α .

Make sure gradient descent is working correctly

Try to make the plot b/w No of iterations



For 100 iteration we will get the value of θ & will evalut

$J(\theta)$ at that point & the height is $J(\theta)$

The plot is showing as No of iterations
↑se the value of $J(\theta)$ is ↓ing

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

Summary

- If α is too small: slow convergence
- If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge

To choose α , try

... 0.001, 0.01, 0.1, 1, ...

Logistic Regression

↳ Classification

↳ $y \rightarrow \text{Discrete O/P}$

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

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

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

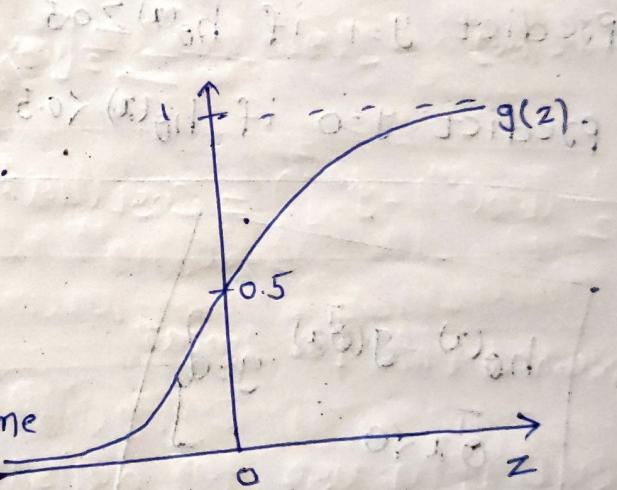
• Logistic Regression [$0 \leq h_{\theta}(x) \leq 1$]

Hypothesis Representation

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

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

↳ Sigmoid fnc
↳ logistic func



Interpretation of hypothesis O/P

$h_{\theta}(x)$ = estimated probability that $y=1$ on iIP x_i

$$\text{eg:- } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

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

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

"Probability that $y=1$, given x , parameterized by θ "

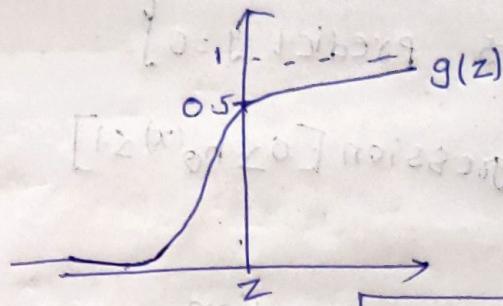
$$P(y=0|x; \theta) = 1 - P(y=1|x; \theta)$$

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

Decision Boundary

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

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



Predict $y=1$ if $h_\theta(x) \geq 0.5$

Predict $y=0$ if $h_\theta(x) < 0.5$

$g(z) \geq 0.5$ when $z \geq 0$

~~get~~
 $y = 1$ when
 $\theta^T x \geq 0$

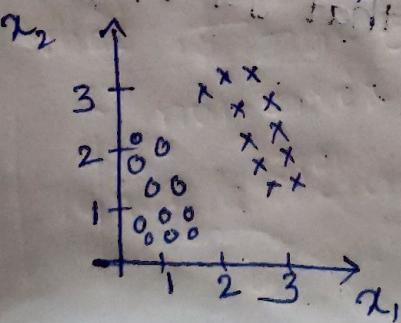
$$h_\theta(x) = g(\theta^T x)$$

$y = d$

$\theta^T x \leq 0$

$x^T \leq 0$

Decision boundary illustration



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

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

$$y = 1 ?$$

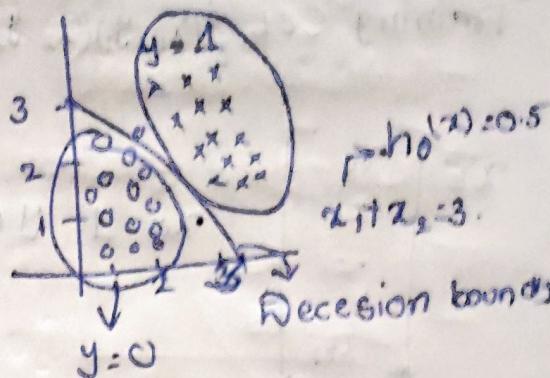
$$y = 0 ?$$

Predict $y=1$ if $z \geq 0$ i.e. $\theta^T x \geq 0$ i.e.

$$-3 + \theta_1 x_1 + \theta_2 x_2 \geq 0 \rightarrow y=1$$

$$-3 + \theta_1 x_1 + \theta_2 x_2 = 0$$

$$\boxed{\theta_1 x_1 + \theta_2 x_2 = 3}$$

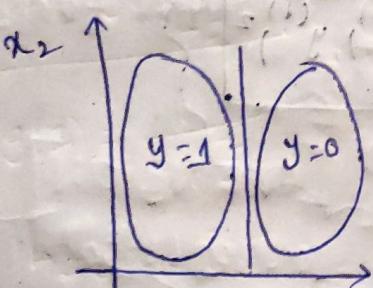


e.g.- Features x_1 & x_2 & $\theta_0 = 5$, $\theta_1 = -1$, $\theta_2 = 0$, So

$$h_0(x) = g(5 - x)$$

Decision boundary of $h_0(x)$?

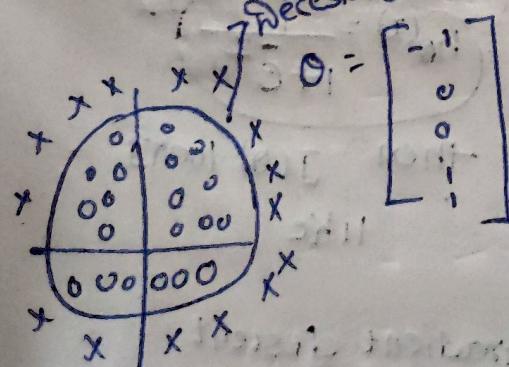
\Rightarrow



$$y=0; 5 - x < 0$$

\Rightarrow Non-linear decision boundaries.

$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$



$$\therefore z = -1 + x_1 + x_2$$

$$\therefore y=1; x_1 + x_2 \geq 1$$

$$y=0; x_1 + x_2 < 1$$

logistic regression cost function

Training Set $S(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_m, y_m)$

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

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

How to choose θ ?

For

Line-Reg

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

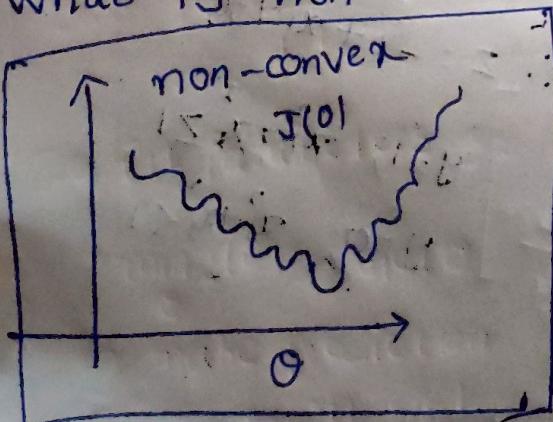
$$\frac{1}{2} [h_{\theta}(x^{(i)}) - y^{(i)}]^2 = \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})^2$$

worked fine for linear

But if u use for logistic then

$[h_{\theta}(x) - y]^2$ will be "non-convex"

what is non-convex?



WKT

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

then $J(\theta)$ looks like

→ Gradient descent will not converge

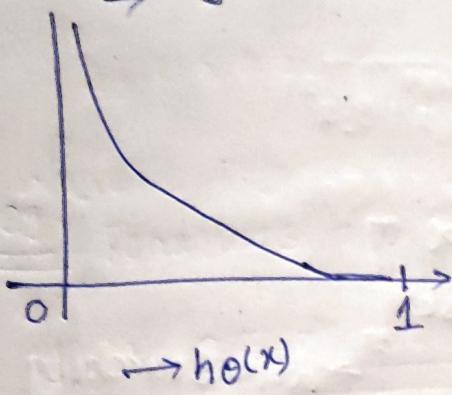
logistic regression Cost function

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

↳ single training example

$\Rightarrow y=1$

In general



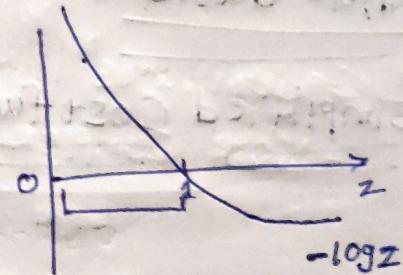
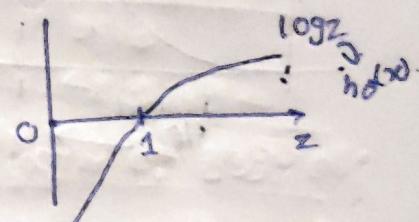
$\rightarrow h_\theta(x)$

Fig

if $y=1$, then $h_\theta(x)=1$

But as $h_\theta(x) \rightarrow 0$

$\rightarrow \text{Cost} \rightarrow \infty$



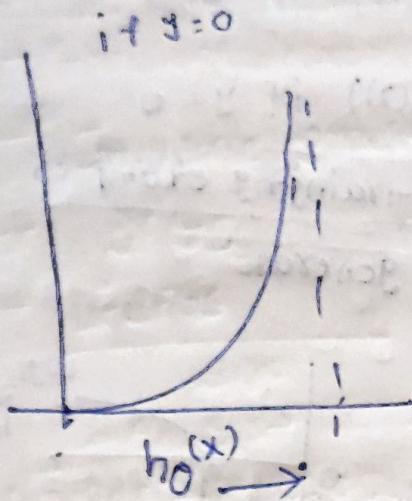
Interested b/w 0 & 1
so, it is as
above in Fig.

Captures intuition that if $h_\theta(x) = 0$,

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

we'll penalize learning algorithm by a very
large cost.

$$\Rightarrow y=0$$



whole dataset

simplified Cost function & Gradient descent

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

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

For $y=0$ or $y=1$

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

To fit parameters θ :

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

To make prediction

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

Gradient Descent

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

Repeat {

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

↑ simultaneously update all θ_j

}

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

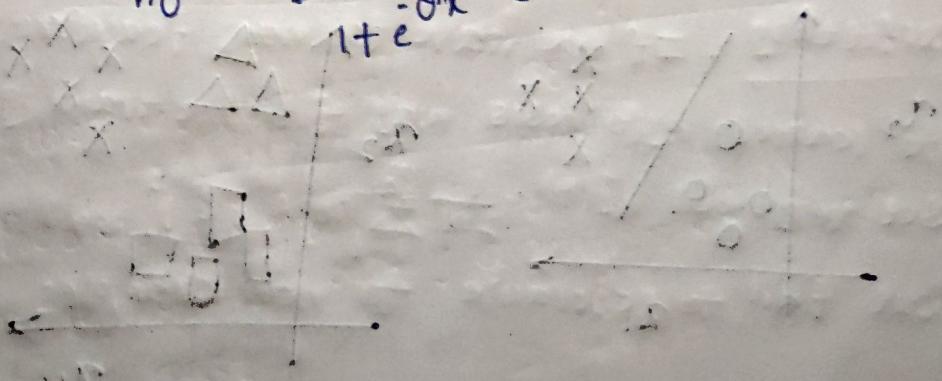
$$\boxed{\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)}}$$

↳ exactly same as linear regression.

but $h_\theta(x) = \theta^\top x$ [for linear regression]

$$h_\theta(x) = \frac{1}{1+e^{-\theta^\top x}} \quad \text{[logistic]}$$

} not same



Advanced optimization

optimization algorithms

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

No need to manually pick α
often faster than Grad-desc

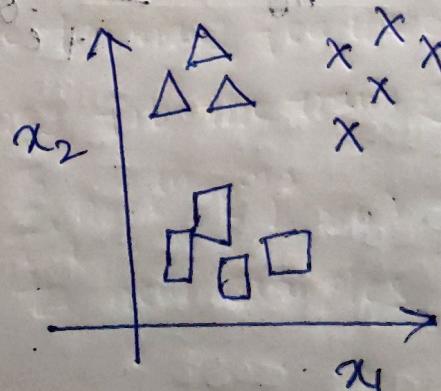
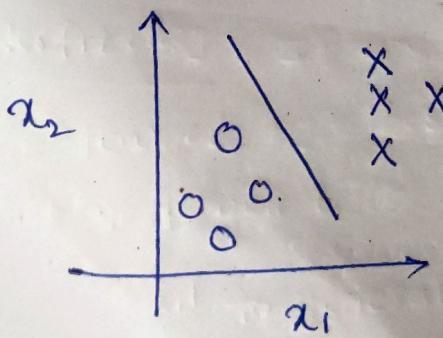
Disadvantages

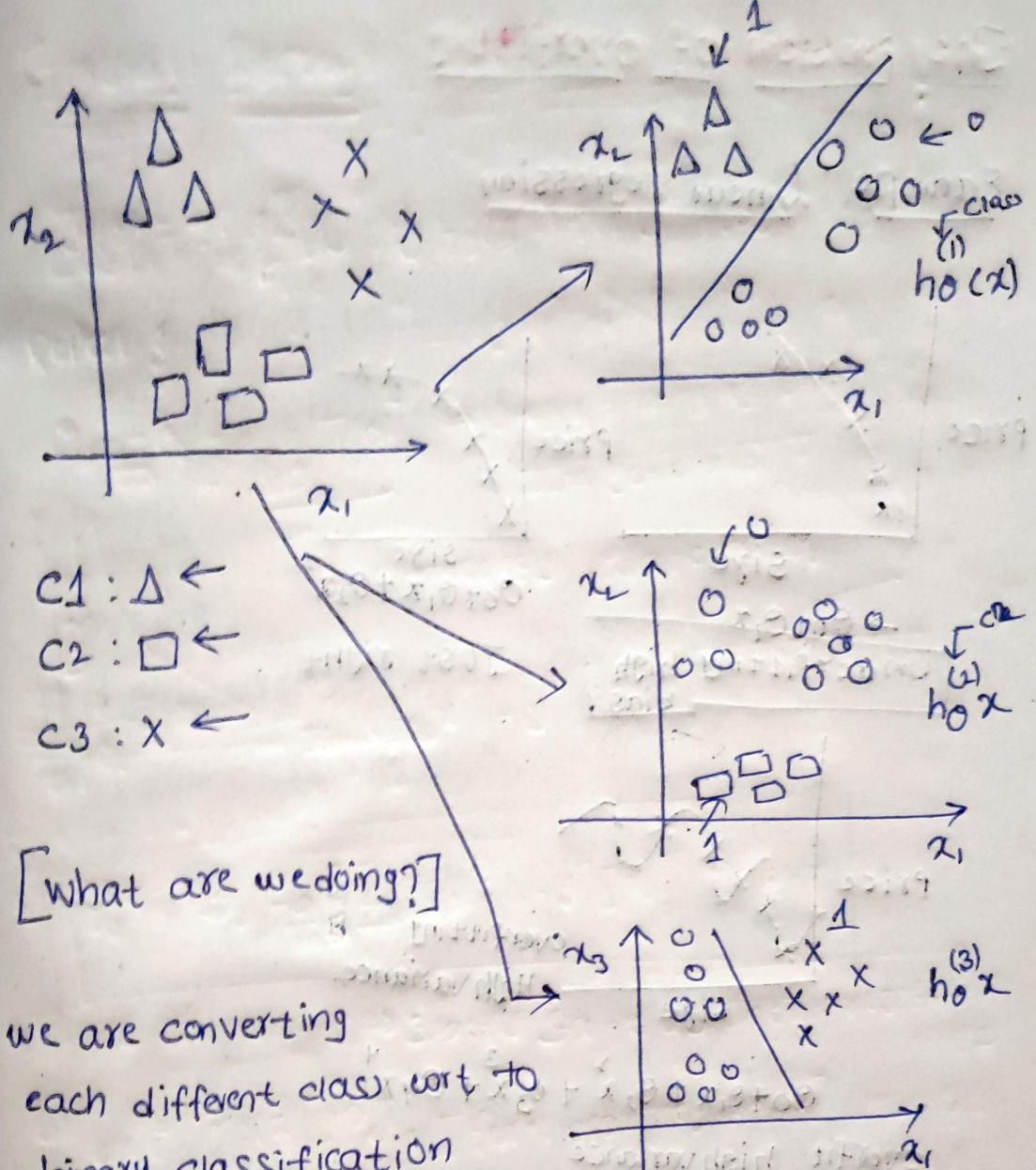
- More Complex

→ line searching
algorithm will
make the
work for α

Multiclass Classification

[One Vs all]





Problem

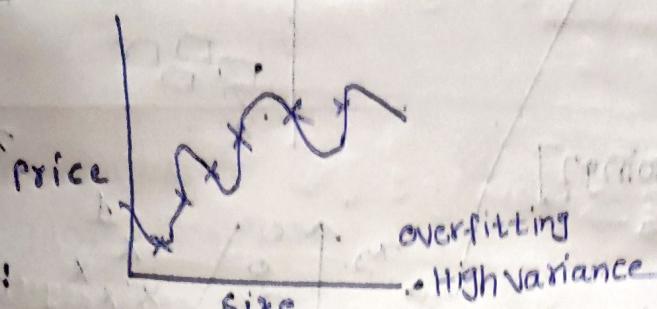
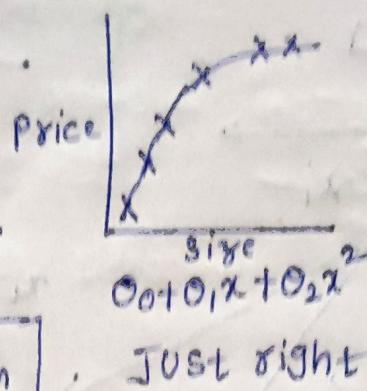
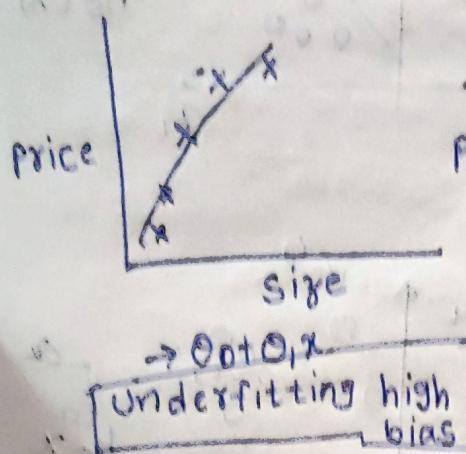
$$h_0^{(i)}(x) = p(y=i|x; \theta) \quad (i=1,2,3)$$

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

$$\underset{i}{\operatorname{max}} h_0^{(i)}(x)$$

The problem of overfitting

Example : Linear regression



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

overfit high variance

Overfitting : If we have too many features.

the learned hypothesis may fit the training

Set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m h_\theta(x^{(i)}) - y^{(i)} \approx 0$),

but fail to generalize to new examples.

Regularization

- Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$
- Less prone to overfitting
- "Simpler" hypothesis

Example

Housing

⇒ Features: x_1, x_2, \dots, x_{100}

⇒ Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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