

- ① Supervising Learning - Logistic Regression, Classification Problem.
- ② Unsupervised Learning: find some Structure (Cluster)

* Linear Regression with One Variable

① Model And Cost Function:

$m \rightarrow$ No. of training examples.

x 's \rightarrow input variable / feature

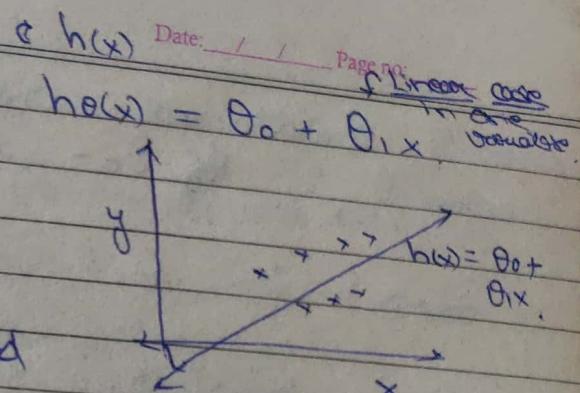
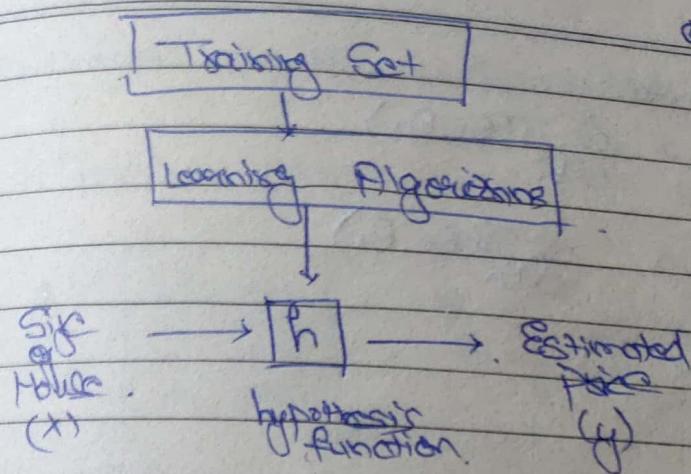
y 's \rightarrow output variable / ~~feature~~ target

$(x, y) \rightarrow$ one training example

$(x_i^0, y_i^0) \rightarrow$ i^{th} training ex.

Training Set

Size in feet (x 's)	Price (y 's)
2104	460
1416	232
x^2	y^2
1	1
1	1
M210	



θ_i 's \rightarrow Parameters

Goal

$$\underset{\theta_0, \theta_1}{\text{Minimize}} \frac{1}{2m} \sum_{i=1}^m [h(x^i) - (y^i)]^2$$

cost function or
sq. error function

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

$$\cancel{J(\theta_0, \theta_1)} = \underset{\theta_0}{\text{Min}} \quad \cancel{\underset{\theta_1}{\text{Min}}}$$

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

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

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

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

- ① Gradient Parameter Learning -
- gradient descent to minimization

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

(2x-1)
①

Date: / / Page no: _____

Correct Simultaneous Update

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

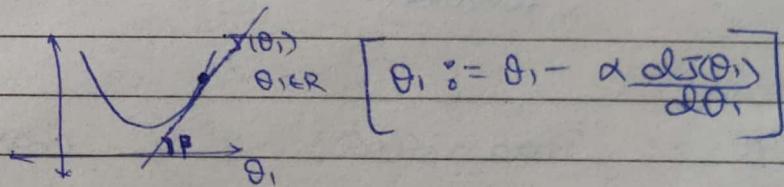
$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$

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

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

Let

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



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

$$\theta_0 \quad (j=0) := \theta_0 - \alpha \cdot \text{grad}$$

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

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

Batch Gradient Descent \rightarrow all train set in each set

③ Linear Algebra Review :-

Date: / / Page no. _____

$$h_{\theta}(x) = -40 + 36x$$

$$x = 2104, 1416, 1534, 852$$

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

$$\theta = \begin{bmatrix} -40 \\ 36 \end{bmatrix}$$

$$C = A \times \theta$$

$$y = ?$$

Prediction = Data Matrix \times Parameter

WEEK - 9

① Multivariate Linear Regression $\begin{array}{c|c|c|c|c} x_1 & x_2 & x_3 & x_4 & y \\ \hline - & - & - & - & \text{---} \end{array}$

n → No. of features.

Vector $\rightarrow x_i^o$ → input (features) of i^{th} example.

x_j^o → Value of feature j in i^{th} example.

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

for convenience $x_0 = 1$

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

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

• vec gradient descent for Multivariate Linear regression.

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

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

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

→ Then

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

• Scaling the feature

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

$$x_2 = \text{No. of bedrooms (1 - 5)}$$

(Free scaling).

$$x_1 = \frac{x_1 - \text{mean}}{\text{range}} \quad x_2 = \frac{\text{No. of bedrooms}}{5}$$

→ Take Mean Normalization only.

$$x_i \rightarrow \frac{x_i - \text{mean}}{\text{range}}$$

$$x_1 = \frac{x_1 - \text{mean}}{\text{range}}$$

Standard deviation

• Learning rate Alpha +

• features & Polynomial Regression +

② Computing Parameters Analytically -

• Normal eq. method to get optimal θ .

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = 0 \rightarrow (\text{for every } j)$$

$$X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & x_3^1 & x_4^1 \\ 1 & \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^n & x_2^n & x_3^n & x_4^n \end{bmatrix}_{m \times (n+1)}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1}$$

$$\boxed{\theta = (X^T X)^{-1} (X^T)(y)} \rightarrow \text{optimal answer.}$$

Scaling not necessary.

WEEK - 3.

① Classification and Representation.

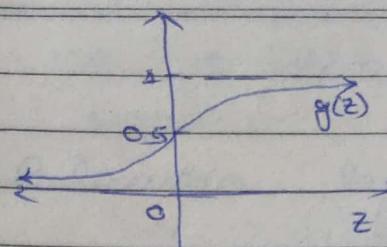
* Logistic Regression - $0 \leq h_\theta(x) \leq 1$

$y \in \{0, 1\} \rightarrow \text{binary class.}$

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

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

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



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

↑ feature.

↳ Probability $y=1$ for given z ,
parametrized by θ .

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

$$\begin{cases} z > 0 \\ g(z) > 0.5 \end{cases}$$

$$\theta^T z > 0$$

Predict $y=1$

$$\theta^T z < 0$$

Predict $y=0$

← Decision boundary.

② Logistic Regression Model

• Cost function:

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \subset \mathbb{R}^{n+1}$$

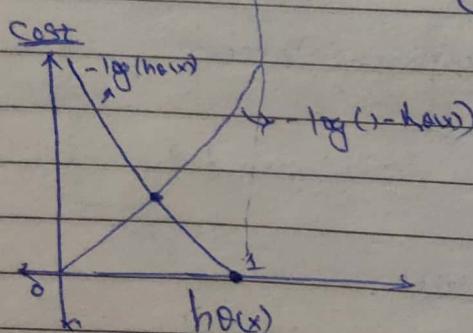
$$x_0 = 1 \quad y \in \{0, 1\}$$

Linear
Regression
cost func.

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

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

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



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

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

- conjugate gradient
- BFGS
- L-BFGS

Date: / / log principle of Maximum likelihood.
Page no. _____

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

$\min_{\theta} J(\theta)$ ← Use Gradient descent

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

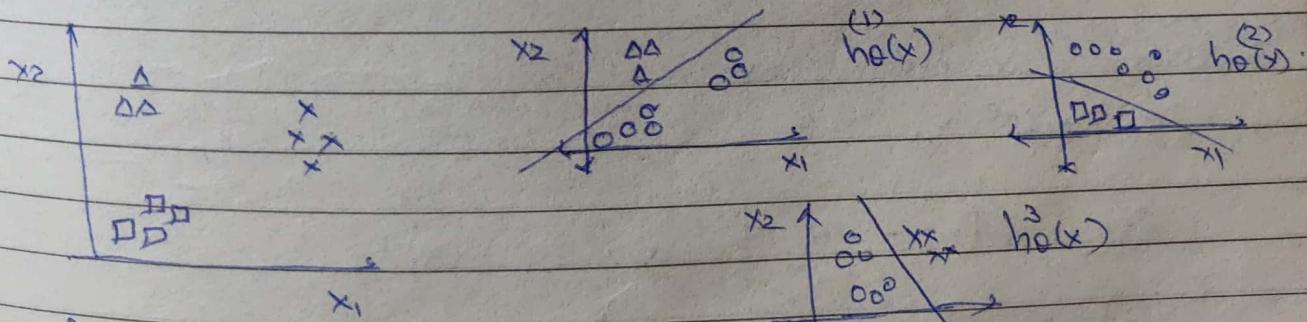
$$\theta_j := \theta_j - \alpha \sum_{i=1}^m [h_\theta(x_i) - y_i] x_i$$

Gradient descent for Logistic regression

Vectorized implementation

$$\theta = \theta - \alpha \frac{1}{m} X^T (g(X\theta) - y)$$

③ Multiclass Classification



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

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

$$X = [2 \times 2]$$

$$y = [2 \times 1]$$

$$\theta = [2 \times 1]$$

$$g(\theta) = [2 \times 1]$$

$$h(x) = \frac{1}{2}x^T\theta$$

$$Date: _____$$

$$Page no.: _____$$

- ④ Solving the problem of
High variance \rightarrow overfitting.
High bias \rightarrow underfitting.

$^{12 \times 1}$
 $^{2 \times 1}$

(2x1)

\rightarrow Solving overfitting

(i) Reduce No of features.

\rightarrow Manually \rightarrow Model Selection Algo (later).

(ii) Regularization.

\rightarrow Keep all features but reduce their magnitude.



• Regularization.

feature $\rightarrow x_1, x_2, \dots, x_{100}$

parameter $\rightarrow \theta_0, \theta_1, \dots, \theta_{100}$

for linear regularization.

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

↑ regularization.
↑ parameter

Gradient descent after
Regularization

$$\theta_0 := \theta_0 - \alpha \sum_{i=1}^m (h_\theta(x_i) - y_i) x_i$$

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

Regularization for logistic regression.

Date: _____
Page no.: _____

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

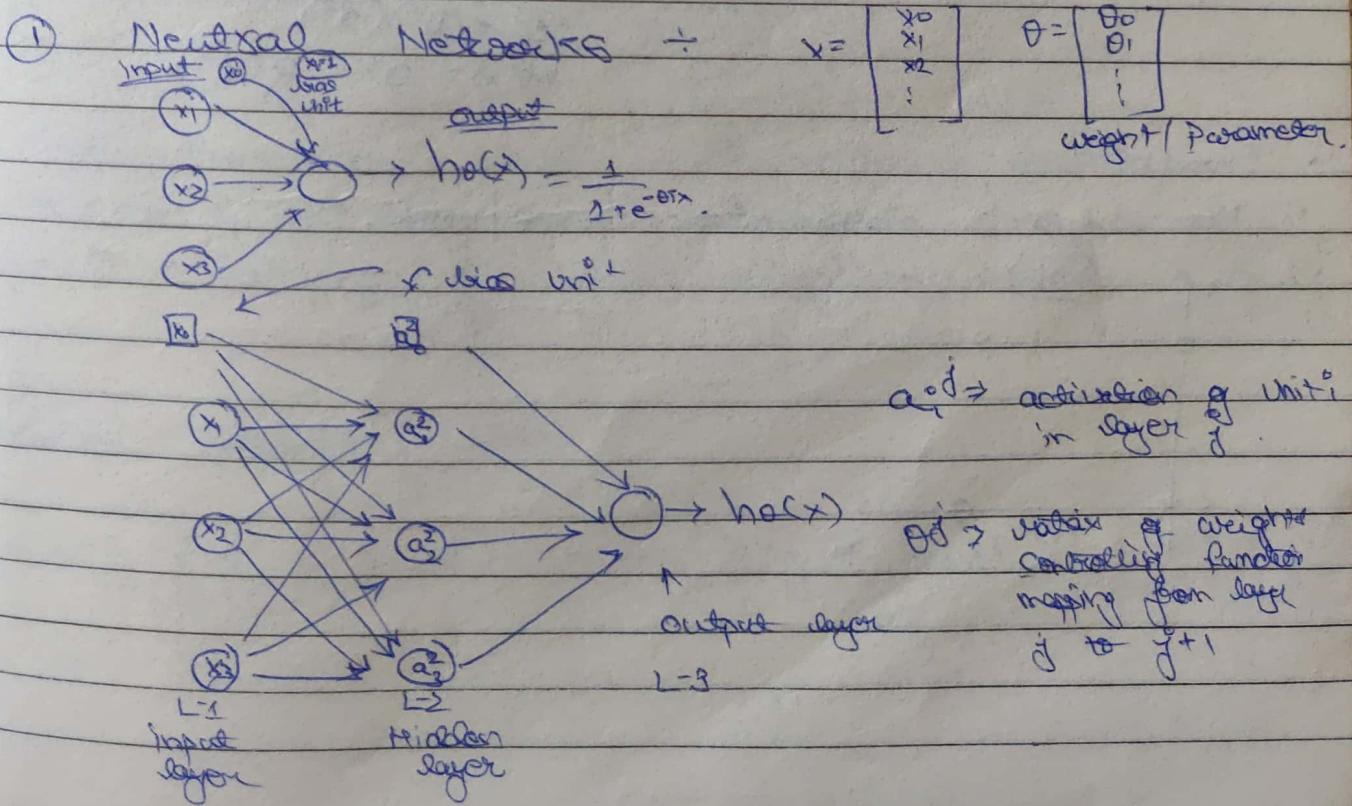
Gradient descent

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

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

WEEK-4.

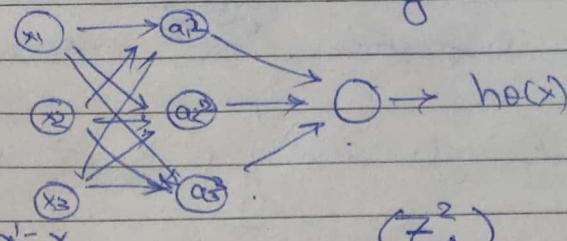
- Non-linear hypothesis need more than features we have to compute $O(n^2)$ total possibilities.



a_0^2

θ_0 θ_1 a_0^2

Date: / / Page no: _____



$\theta' \in \mathbb{R}^{3 \times 4}$

$x = [x_0, x_1, x_2, x_3]$

$a^2 = x$

$$a_1^2 = g(\theta_{10} x_0 + \theta_{11} x_1 + \theta_{12} x_2 + \theta_{13} x_3)$$

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

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

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$z^2 = \begin{bmatrix} z_1^2 \\ z_2^2 \\ z_3^2 \end{bmatrix}$$

$$z^2 = (\theta' x) \alpha (\theta' a^2)$$

$$a^2 = g(z^2) \quad R^3$$

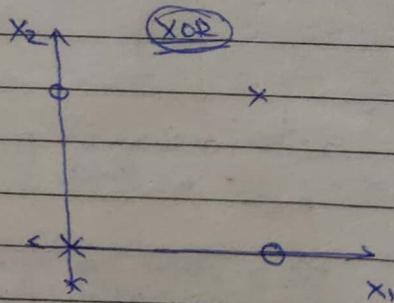
$$A_{act} \quad \star \star \quad a_0^2 = 1 \quad R^4$$

$$z^3 = (\theta^2)(a^2)$$

$$h_\theta(x) = a^3 = g(z^3)$$

• Examples:

• Non-Linear classification XOR/XNOR.



① C

SU

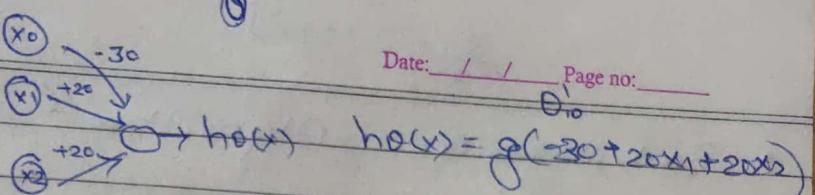
Cost

Neut

J(e) =

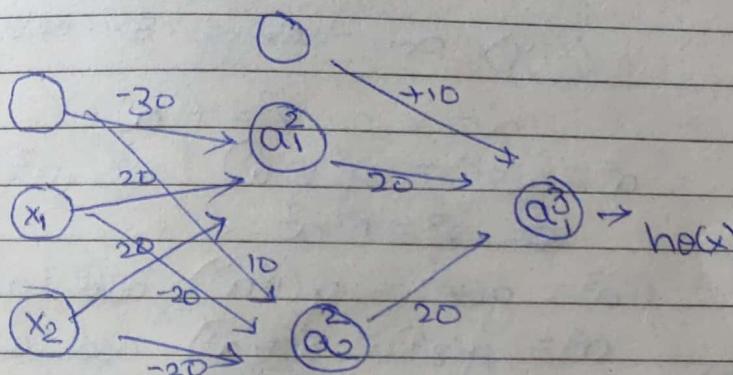
$$x_1, x_2 \in \{0, 1\}$$

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



x_1	x_2	$h(x)$
0	0	$g(-30) = 0$
1	0	$g(-10) = 0$
1	1	$g(10) \approx 1$
0	1	$g(-10) = 0$

$$\text{XOR} = (\text{A AND B}) \cdot (\bar{\text{A}} \text{ AND } \bar{\text{B}})$$



x_1	x_2	a_1^2	a_2^2	$h(x)$
0	0	0	1	0
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

WEEK 5

① Cost function And Back propagation

s_l = No. of unit (Not biased) in layer l .

$$\text{Cost fn} \Rightarrow J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m y_i^* \log h(x_i^*) + (1-y_i^*) \log (1-h(x_i^*)) \right] + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

Neural Network $R^{L(K)}$

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^* \log h(x_i^*)_k + (1-y_k^*) \log (1-h(x_i^*)_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{s=1}^{S_l} \sum_{j=1}^{S_{l+1}} \left[\theta_j^* \right]^2$$

$$\frac{3(1.01)^4 + 4 - [3(0.99)^4 + 4]}{(0.02)}$$

$$\frac{3[(1.01)^4 - (0.99)^4]}{(0.02)}$$

Date: / / Page no: _____

• Backpropagation Algorithm :-

$$\text{Min } J(\theta)$$

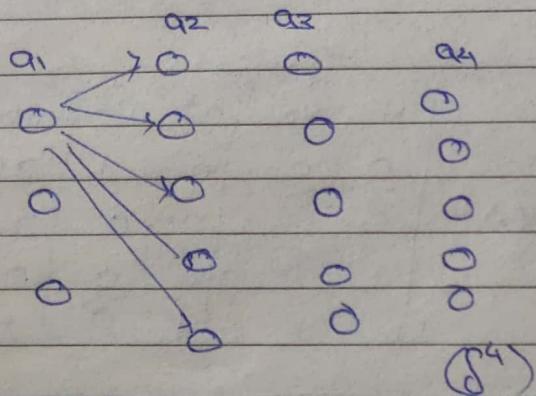
$$\boxed{\begin{array}{c} J(\theta) \\ \frac{\partial J(\theta)}{\partial \theta_j} \end{array}}$$

**

②

.

Forward Propagation.



(x, y) one testing Example

$$a^1 = x \quad z^2 = \theta^1 a^1$$

$$a^2 = g(z^2) = g(\theta^1 a^1) \text{ add } (a_0^2)$$

$$z^3 = \theta^2 a^2$$

$$a^3 = g(z^3) = g(\theta^2 a^2) \text{ add } (a_0^3)$$

$$a^4 = g(z^4) = g(\theta^3 a^3) \text{ add } (a_0^4)$$

X

x =

$\theta_1 =$

Backpropagation Algorithm

$s_j^l = \text{error}$ if node j in layer l .

$$s_j^{l-1} = (a_j^{l-1})^T - y_j \quad s^4 = a^4 - y$$

$$s^3 = (\theta^3)^T s^4 * \overbrace{g'(z^3)}^{a^3(1-a^3)}$$

$$s^2 = (\theta^2)^T s^3 * g'(z^2)$$

No s^1 back no error.

$$\begin{array}{r} 15 \\ \times 24 \\ \hline 39 \end{array}$$

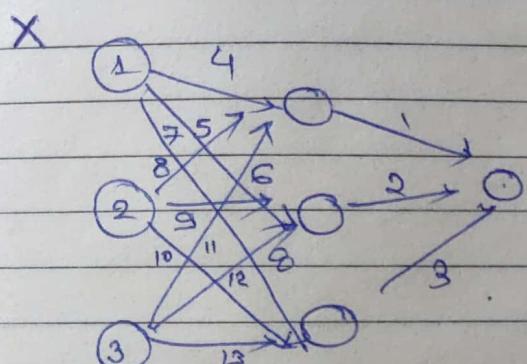
Date: _____ / _____ / _____ Page no: _____

$$\frac{d}{d\theta y^2} J(\theta) = \alpha y^2 \delta_i^{(x+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda=0)$$

* Complete Backpropagation Watch ScreenShot *

② Back propagation in detail:

- # gradient checking



$$a_s^2 = g(4)$$

$$\theta^2 = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$5000 \times 25$$

5003 x 25

$$x = \boxed{\quad}$$

5000x400

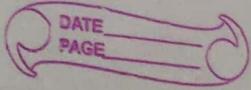
$$J = \begin{bmatrix} & \\ & \end{bmatrix}_{5000 \times 1}$$

$$\theta_1 = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$B_2 = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$h(\theta) = \boxed{[]}$$

WEEK - 6



- What should we do?
- Get More training examples → fix high variance check
- Try smaller set of features → fix high variance
- Try getting additional features → fix high bias
- Try adding polynomial feature → fix high bias
- Try decreasing → fix high bias
- Try increasing → fix high variance

- Evaluating a Hypothesis:

divide Training Set

2 Parts

$$\frac{[]}{[]} \xrightarrow{\text{Training Set 70\%}} \xrightarrow{\text{Testing Set 30\%}} \xrightarrow{\text{randomly}}$$

Computing Test error

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^i) - y_{\text{test}}^i)^2$$

for logistic Test error:

$$J_{\text{test}}(\theta) = \frac{-1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left[y_{\text{test}}^i \log h_{\theta}(x_{\text{test}}^i) + (1-y_{\text{test}}^i) \log (1-h_{\theta}(x_{\text{test}}^i)) \right]$$

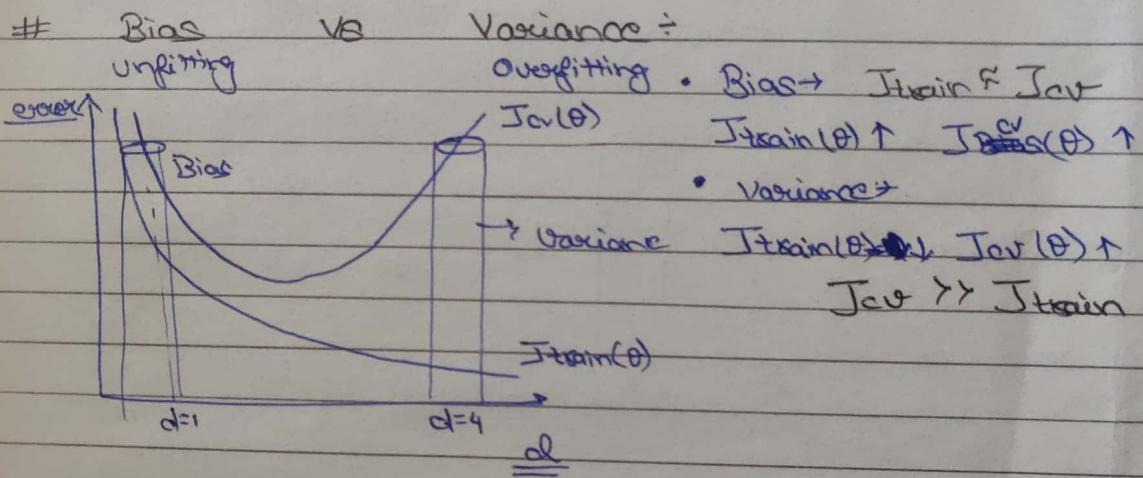
for classification - Min classification.

$$err(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \leq 0.5 \text{ & } y = 0, \text{ or, } h_{\theta}(x) \geq 0.5 \text{ & } y = 1 \\ 0 & \text{otherwise.} \end{cases}$$

$$\text{Test Error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} err(h_{\theta}(x_{\text{test}}^i), y_{\text{test}}^i)$$

• Model Selection and Training Set.

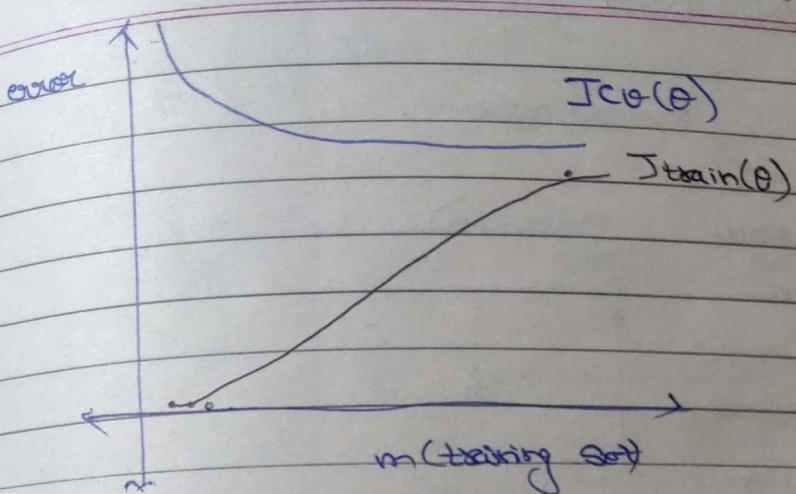
- Breaking of Training Set → Optimize the θ
- ① Training Set (80%) → final Polynomial
 - ② Cross Validation Set (20%) → degree d
 - ③ Test Set (20%) → $J_{\text{test}}(\theta)$



• Learning Curves

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

$$J_{\text{cv}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}_{\text{cv}}) - y^{(i)}_{\text{cv}})^2$$



Machine Learning Design System :-

① Build a spam classifier :-

② Handling skewed data :-

		Actual class	
		1	0
Predicted class	1	True +ve True -ve	False +ve False -ve
	0	False -ve	True -ve

$$\text{Precision} = \frac{\text{True +ve}}{\text{True +ve} + \text{False +ve}}$$

$$\text{Recall} = \frac{\text{True +ve}}{\text{True +ve} + \text{False -ve}}$$

Trending off b/w precision and recall
use F-Scores.

$$\frac{2ab}{a+b}$$

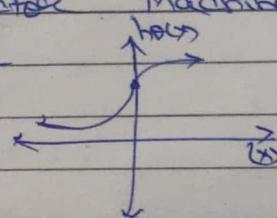
WEEK-7

Large Margin Classification

- SVM (Support Vector Machine)

$$z = \theta^T x$$

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



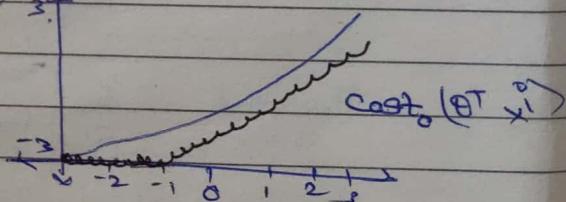
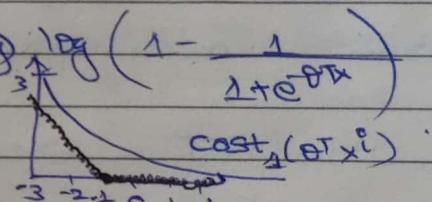
if $y=1$ $[\theta^T x > 0]$
 $y=0$ $[\theta^T x \leq 0]$

from Logistic Regression

$$J_{\theta} = -y \log \frac{1}{1 + e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}} \right)$$

$$y=1 \quad J_{\theta} = -\log \left(\frac{1}{1 + e^{-\theta^T x}} \right)$$

$$y=0 \quad J_{\theta} = -\log \left(1 - \frac{1}{1 + e^{-\theta^T x}} \right)$$

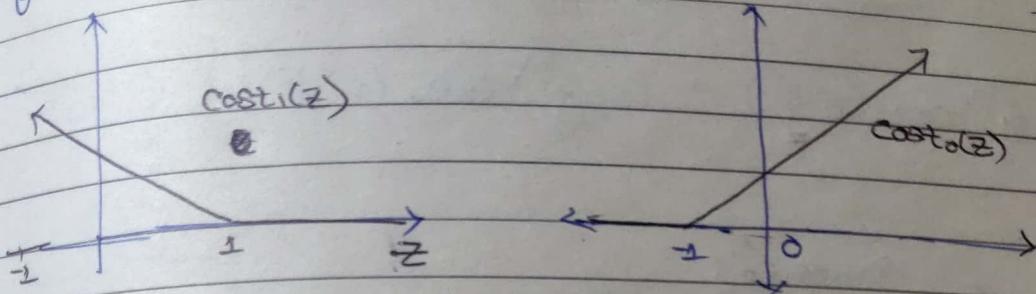


SVM Cost function :

$$J(\theta) = \min_{\theta} C \left[\sum_{i=1}^m y^i \text{cost}(\theta^T x^i) + (1-y^i) \text{cost}(\theta^T x^i) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$h_\theta(x) = \begin{cases} 1 & \text{if } \theta^T x > 0 \\ 0 & \text{if } \theta^T x \leq 0 \end{cases}$$

$$\underset{\theta}{\text{Min}} \sum_{i=1}^m \left[y_i \text{Cost}_+ (\theta^T x_i) + (1-y_i) \text{Cost}_0 (\theta^T x_i) \right] + \frac{1}{2} \sum_{j=1}^q \theta_j^2$$



if $y=1$ $\hat{y} \geq 1$ ($\theta^T x \geq 1$)
 $y=0$ $\hat{y} \leq -1$ ($\theta^T x \leq -1$)

Kernels :

$f_1 = x_1 \quad f_2 = x_2 \quad f_3 = x_1 x_2 \quad \dots$ (features)
 ↳ To Think better choice of features.

→ Gaussian Kernel.

$$f_i = \text{Similarity}(x, l^{(i)}) = \exp \left[- \frac{\|x - l^{(i)}\|^2}{2\sigma^2} \right] = \exp \left[- \frac{\sum_j (x_j - l_j^{(i)})^2}{2\sigma^2} \right]$$

$$f_i = K(x, l^{(i)})$$

$$\textcircled{1} \quad x \approx l^1$$

$$\hookrightarrow f_1 \approx 1$$

$$\textcircled{2} \quad x \text{ is far from } l^1$$

$$\hookrightarrow f_1 \approx 0$$

How to get landmark.

$$\underset{\theta}{\text{Min}} \ C \sum_{i=1}^m (y_i) \text{cost}_1(\theta^T f_i) + (1-y_i) \text{cost}_0(\theta^T f_i) + \frac{1}{2} \sum_{j=1}^m \theta_j^2.$$

SVM in Practice :-

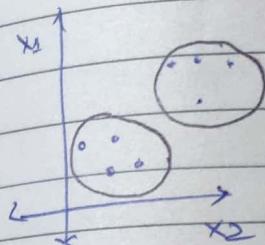
liblinear, libsvm

- ① Choice of Parameter C
- ② Choice of Kernel.

• Objec

C
y

Clustering (UnSupervising learning)



K-Means Algorithm :-

- ① cluster assignment
- ② Move centroid.

• Objective function of K-mean.

 $c^i \Rightarrow (1 \times n_k)$ No. of cluster. $y_k \Rightarrow$ cluster centroid. ($y_k \in R^n$) $y_{ci} \Rightarrow x^i \in S \quad y_{ci} = y_5$
 $c^i = 5$ • Objective function \Rightarrow

$$J(c^1, \dots, c^m; y_1, \dots, y_k) = \frac{1}{m} \sum_{i=1}^m \|x^i - y_{ci}\|^2$$

$$\boxed{\min_{c^1, \dots, c^m, y_1, \dots, y_k} J(c^1, \dots, c^m, y_1, \dots, y_k)}$$

Fahad Shrikh

Dimensionality Reduction (unsupervised).

(more feature \rightarrow less features.

$$3D \leftrightarrow 2D \leftrightarrow 1D$$

Principle Component Analysis:

[PCA] (reduce dimension).

Reduce n vector to k-Vector. Project

K vector u^1, u^2, \dots, u^K to Minimize error.

① Preprocessing (mean normalization & scaling)

② n-dimension to k-dimension

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (x^i) (x^i)^T \quad (n \times n) \text{ matrix}$$

compute eigenvector

$$*** [U, S, V] = \text{SVD}(\Sigma);$$

$$U \rightarrow [n \times n] \rightarrow \begin{bmatrix} | & | & \dots & | \\ u^1 & u^2 & \dots & u^n \\ | & | & \dots & | \end{bmatrix}$$

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$z = \begin{bmatrix} | & | & \dots & | \\ u^1 & u^2 & \dots & u^k \\ | & | & \dots & | \end{bmatrix} \quad n \times k$$

$$U_{\text{reduce}} = U(:, 1:k)$$

$$Z = U_{\text{reduce}}^{-1} \times x^i$$

$$x = (Z^T) (Z) \Rightarrow k \times 1$$

DATE
PAGE

Reconstruction from compressed representation :-

$$\text{Yapp} = \text{Usedur} \circ Z^{(1)} \\ R^n \qquad \qquad \qquad \qquad \qquad K_1$$

WEEK - 9

DATE _____
PAGE _____

Density Estimation :-

↳ density at outside stand alone

↳ Anomaly detection.

$$P(x) < \epsilon$$

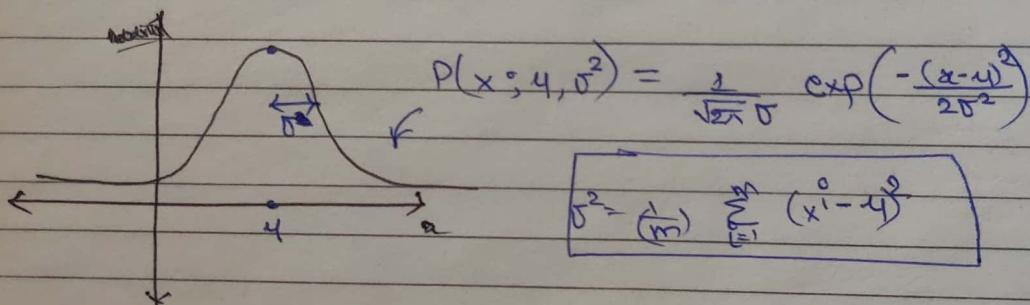
Gaussian (Normal) Distribution :-

Say $x \in \mathbb{R}$. If x is a distributed Gaussian with Mean μ , Variance σ^2 ,

it distributed as

$$x \sim N(\mu, \sigma^2)$$

$$\mu = \left(\frac{1}{m} \right) \sum_{i=1}^m (x_i)$$



Anomaly Detection Algorithm :-

(for very few the example) (Not use Supervising Learning algorithm)

$$P(x_1; \mu_1, \sigma_1^2) \times P(x_2; \mu_2, \sigma_2^2) \times \dots \times P(x_n; \mu_n, \sigma_n^2)$$

$$\Rightarrow \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$

- Anomaly detection
- Fraud detection
- Manufacturing (aircrafts)
- Monitoring Machine

- Supervised Learning
- Email spam
- Weather prediction
- Cancer classification

Multivariate Gaussian Distribution:

①

$$P(x^0, y, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-y)^T \Sigma^{-1} (x-y)\right)$$

$$y \in \mathbb{R}^m \quad \Sigma \in \mathbb{R}^{n \times n}$$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^i$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^i - \mu)(x^i - \mu)^T$$

② Now our new example
($\hat{x} \rightarrow \mu$)

* Multivariate Gaussian will be original Gaussian variation if their contours plot on axis aligned mean

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}$$

Recommendation Systems :-

$n_u \Rightarrow$ No. of users

$n_m \Rightarrow$ No. of movies.

$x(i,j) \Rightarrow 1$ if user j rated movie i .

$(y^{i,j}) \Rightarrow$ rating user j to movie i (if $x(i,j)=1$).

Movie	A	B	C	D	x_1 (Param)	x_2 (Const)
1 $x_1 \mapsto$	5	5	0	0	0.9	0
2 $x_2 \mapsto$	5	4.5	0	0	1.0	0.01
3	0.5	4	0	0	0.99	0
4	0	0	5	4	0.1	1.0
5	0	0	5	4	0	0.9

① Content based recommendation :-

x_1, x_2 feature.
 $n=2$ (feature)

$$x^i = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix} \rightarrow \text{intercept}$$

* for each user j , learn parameter $\theta^j \in \mathbb{R}^{n+1}$
user j as rating movie i with $(\theta^j)^T (x^i)$ stars.

$$\theta^j \quad x^3 = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \quad \theta^j = \begin{bmatrix} 0 \\ 0.99 \\ 0 \end{bmatrix} \quad (\theta^j)^T (x^3) = (5)(0.99) = 4.95.$$

* To learn θ^j (use linear regression) :-

$$\text{Min } \frac{1}{2m} \sum_{i: x(i,j)=1} \left[(\theta^j)^T (x^i) - y^{i,j} \right]^2 + \frac{\lambda}{2m} \sum_{k=1}^n (\theta_k^j)^2$$

$\theta^j \in \mathbb{R}^{n+1}$

Formally
To learn (θ^i)

$$\text{Min}_{(\theta^i)} \left(\frac{1}{2} \right) \sum_{i: x_i, j=1} \left[(\theta^i)^T (x^i) - (y^{i,j}) \right]^2 + \lambda \sum_{k=1}^n (\theta_k^i)^2$$

To learn $\theta^1, \theta^2, \dots, \theta^{nm}$

$$\text{Min}_{\theta^1, \dots, \theta^{nm}} \left(\frac{1}{2} \right) \sum_{j=1}^{nm} \sum_{i: x_i, j=1} \left[(\theta^i)^T (x^i) - (y^{i,j}) \right]^2 + \lambda \sum_{j=1}^{nm} \sum_{k=1}^n (\theta_k^j)^2$$

$$J(\theta^1, \dots, \theta^{nm}).$$

To Gradient descent update:

$$(\theta_k)^j = \theta_k^j - \alpha \sum_{i: x_i, j=1} \left[(\theta^i)^T x^i - (y^{i,j}) \right] (x_k^i) \quad (k=0)$$

$$(\theta_k)^j = \theta_k^j - \alpha \left[\sum_{i: x_i, j=1} \left[(\theta^i)^T (x^i) - (y^{i,j}) \right] (x_k^i) + \lambda \theta_k^j \right] \quad (k \neq 0)$$

$\frac{\partial}{\partial \theta_k} J(\theta^1, \dots, \theta^{nm})$.

② Collaborative filtering algorithm →
(automatic feature adding).

Given $(\theta^1, \dots, \theta^{nm})$ to learn (x^i)

$$\frac{1}{2} \sum_{j: x_i, j=1} ((\theta^i)^T (x^i) - (y^{i,j}))^2 + \frac{1}{2} \sum_{k=1}^n (x_k^i)^2.$$

Given $(\theta^0, \dots, \theta^{n_m})$ to learn x^1, \dots, x^{n_m} .

$$\text{Min}_{(x^1, \dots, x^{n_m})} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:y^i = \delta(i,j)=1} \left[(\theta^j)^T (x^i) - (y^i)^j \right]^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^i)^2$$

Gradient descent think yourself.

When θ^j use algo.

$$\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow x$$

Guessing first.

(OR)

Combining

where

the con+

$$J(x^1, \dots, x^{n_m}, \theta^0, \dots, \theta^{n_m}) = \frac{1}{2} \sum_{(y^i): \delta(i,y^i)=1} \left[(\theta^j)^T (x^i) - (y^i)^j \right]^2 + \left(\frac{\lambda}{2} \right) \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^i)^2 + \left(\frac{\lambda}{2} \right) \sum_{j=1}^{n_m} \sum_{k=1}^n (\theta_k^j)^2$$

$$\cancel{\text{Min}}_{(x^1, \dots, x^{n_m})} \quad J(x^1, \dots, x^{n_m}, \theta^0, \dots, \theta^{n_m}) \quad \theta \rightarrow x \rightarrow \theta$$

① Initialize $(x^1, \dots, x^{n_m}) (\theta^0, \dots, \theta^{n_m})$ to small random value

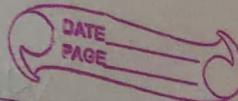
② Min $J(x^1, \dots, x^{n_m}, \theta^0, \dots, \theta^{n_m})$

$$x_{k^i} = x_{k^i} - \alpha \left[\sum_{j: \delta(i,j)=1} ((\theta^j)^T (x^i) - (y^i)^j) \theta_{k^i} + \lambda x_{k^i} \right]$$

$$\theta_{k^i} = \theta_{k^i} - \alpha \left[\sum_{j: \delta(i,j)=1} \left[(\theta^j)^T x^i - (y^i)^j \right] x_k^i + \lambda \theta_{k^i} \right]$$

Vectorization
(Low-Rank
Matrix)

Implementation:



$$X = \begin{bmatrix} -(x_1)^T \\ -(x_2)^T \\ \vdots \\ -(x_m)^T \end{bmatrix}$$

$$\Theta = \begin{bmatrix} -(\Theta^1)^T \\ -(\Theta^2)^T \\ \vdots \\ -(\Theta^k)^T \end{bmatrix}$$

Predicting Rating = $X \Theta^T$

Mean Normalization:
(for Rating at 0 to 1)

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & 2 & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

$$M = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix}$$

$$Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.25 & 1.25 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

for user j on movie i

Learn (Θ^j, x^i)

$$(\Theta^j)^T (x^i) + (y_i)$$

$$\Theta^j = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

~~(Theta)^T x^i~~

WEEK - 10

DATE _____
PAGE _____

* dealing with larger dataset →

① Plot Learning Curve to check. adding Training set will help or not.

Stochastic Gradient descent :

Batch

$$h(\theta) = \sum_{j=0}^n \theta_j x_j$$

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m (h(\theta)(x_i) - y_i)^2$$

{

$$\theta_j = \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h(\theta)(x_i) - y_i) x_i \quad (\text{for } j=0, \dots, n)$$

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

,

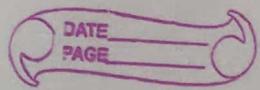
,

,

,

,

,



WEEK-11

Photo OCR :-

