

Statistics : Science of collecting, organizing, analysing and interpreting Data

Data

Nominal

Ordinal

Interval

Variable \rightarrow Quantitative

\rightarrow Discrete

\rightarrow no. of student

\rightarrow Continuous

\rightarrow weight

\rightarrow height

\rightarrow Qualitative

(Categorical)

Ordinal

\rightarrow Feedback

1-10

\rightarrow Rank

Nominal

\rightarrow Eye color

\rightarrow Blood Group

* Machine learning

\Rightarrow Supervised learning

\Rightarrow Unsupervised learning

\Rightarrow Reinforcement learning : Feedback based.

General M.L. Model

I/P \rightarrow Machine \rightarrow Output

Feedback

⇒ M.L. is possible because of large data available
↳ Computation power.

⇒ Application :- Healthcare
:- Sentiment analysis
:- Fraud detection

⇒ M.L. is the science of making computers learn & act like humans by feeding data & information without being explicitly coded / programmed.

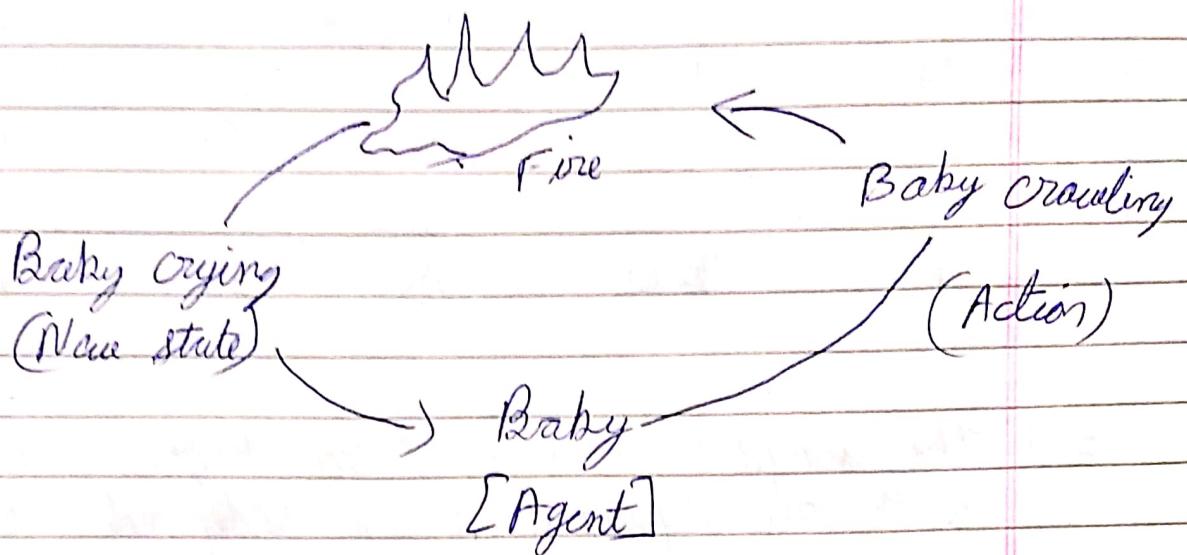
* Steps to M.L. Project [general]

- Define your objective
- Collecting data
- Prepare data
- Select Algorithm
- Train Model
- Test Model
- Predict
- Deploy

* Division based on what M.L. does.

- 1) Want to predict category :- Classification
- 2) Want to ^{predict} quantity :- Regression.
- 3) Detect anomaly :- anomaly detection

* Reinforcement learning is an important type of M.L. where an agent learns how to behave in an environment by performing actions & see the results.



* Supervised & Unsupervised

- Labeled data → Non-labeled data
- Direct feedback → No feedback
- Predict Output → Find hidden structure in data

* Linear Regression

⇒ A statistical model used to predict the relationship between independent & dependent variables

Examine two factors

⇒ Which variables in particular are significant predictors of the outcome variables

⇒ How significant is the regression line to make predictions with highest possible accuracy.

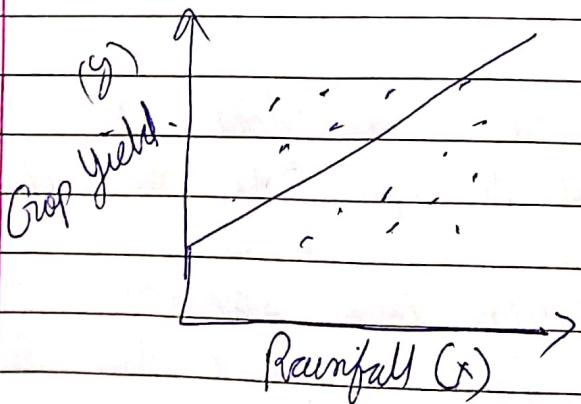
Linear Regression

x	y	(x^2)	(y^2)	$(x \cdot y)$
1	2	1	4	2
2	4	4	16	8
3	5	9	25	15
4	4	16	26	16
5	5	25	25	25
<hr/>				
Σ	15	55	86	66

\Rightarrow The simplest form of linear Regression equation with one dependent & one independent variable is represented by:-

$$y = m \cdot x + c$$

$y \rightarrow$ Dependent variables
 $x \rightarrow$ Independent variables
 $m \rightarrow$ Slope of the line
 $c \rightarrow$ Coefficient of line



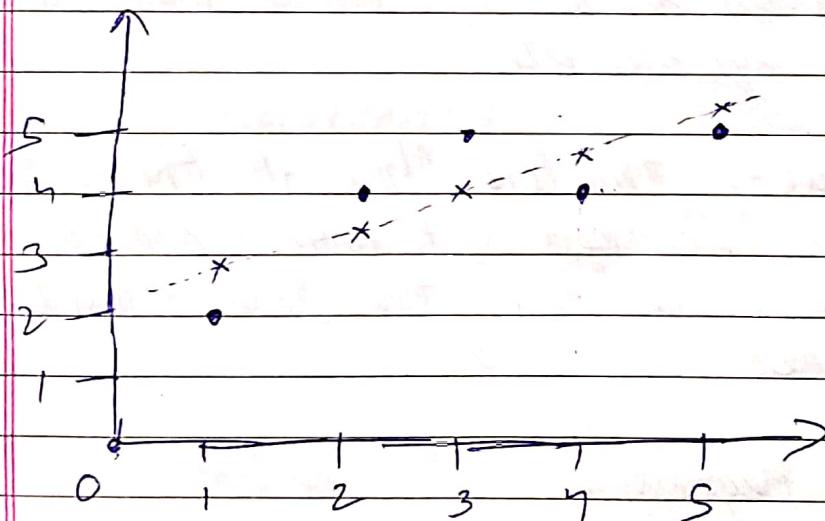
$$m = \frac{((n * \sum(x * y)) - (\sum(x) * \sum(y)))}{((n * \sum(x^2)) - (\sum(x))^2)}$$

$$= \frac{((5 * 66) - (15 * 26))}{((5 * 55) - (125))} = 0.6$$

$$c = \frac{((\sum(y) * \sum(x^2)) - (\sum(x) * \sum(x * y)))}{((n * \sum(x^2)) - (\sum(x))^2)}$$

$$= 2.2$$

$$\begin{aligned}y &= m * x + c \\&= 0.6 * 3 + 2.2 \\&= 4\end{aligned}$$



y_{pred}

$$y = 0.6 * 1 + 2.2 = 2.8$$

$$y = 0.6 * 2 + 2.2 = 3.4$$

$$y = 0.6 * 3 + 2.2 = 4$$

$$y = 0.6 * 4 + 2.2 = 4.6$$

$$y = 0.6 * 5 + 2.2 = 5.2$$

L.R.: Works for continuous values

Linear Supervised learning

Date _____
Page _____

X	\bar{X}	y_{pred}	$(y - y_{pred})$	$(y - y_{pred})^2$
1	2	2.8	-0.8	0.64
2	4	3.4	0.6	0.36
3	5	4	1	1
4	4	4.6	-0.6	0.36
5	5	5.2	-0.2	0.04
Σ				2.4

\Rightarrow The sum of squared errors for this regression line is 2.4. We check this error for each line and conclude the best fit line having the least L square value.

\Rightarrow Minimizing this distance, there are lots of ways, between line and data point like sum of squared errors, sum of absolute errors Root mean square etc.

\Rightarrow We keep moving this line through the data points to make sure the best fit line has the least square distance between the data points & regression line.

M. Multi linear Regression

L.R.

$$y = m * x + c$$

M.LR.

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + \dots + c$$

Slope

Dependent Variable Independent Variables Co-efficient C

Implementation of L.R

- import numpy as np
- import matplotlib.pyplot as plt
- import pandas as pd
- import Seaborn as sns
- matplotlib inline
- # importing the data set & extract dependent and independent variable

```
companies = pd.read_csv('... .csv')
```

```
x = companies.iloc[:, :-1].values
```

```
y = companies.iloc[:, 4].values
```

```
companies.head()
```

- # data visualization, building correlation matrix

```
sns.heatmap(companies.corr())
```

- # encoding categorical data

```
from sklearn.preprocessing import LabelEncoder,  
OneHotEncoder
```

```
labelencoder = LabelEncoder()
```

```
x[:, 3] = labelencoder.fit_transform(x[:, 3])
```

onehotencoder = OneHotEncoder(categorical_features=[3,7])

x = onehotencoder.fit_transform(x).toarray()

print(x[0])

x = x[:, 1:]

- # Splitting of our data from sklearn.model_selection import train_test_split

x-train, x-test, y-train, y-test

= train_test_split(x, y, test_size=0.2, random_state=0)

- from sklearn.linear_model import LinearRegression

regressor = LinearRegression()

regressor.fit(x_train, y_train)

- # Predicting the test set results

y_pred = regressor.predict(x_test)

y_pred =

- # Calculating co-efficient & intercept

`print (regressor.coef_)` # [gives slope]

`print (regressor.intercept_)` # [gives constant]

- # Calculating R squared matrix (evaluating our model)

`from sklearn.metrics import r2_score
r2_score(y-test, y-pred)`

* Logistic Regression

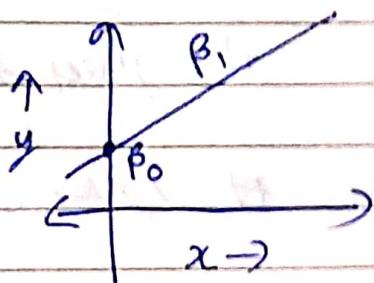
- => An algorithm to perform an binary classification
- => It is of type which gives discrete values
- => supervised learning
- => Odds of success (O) = Probability of an event happening

Probability of an event not happening

$$\text{Or } O = \frac{P}{1-P}$$

range of values of odds is 0 to ∞
values of probability change from 0 to 1

* Maths Behind Logistic Regression



Here β_0 is the y-intercept, β_1 is slope of the line x is value of the x-coordinate y is value of prediction

$$\text{eq: } y = \beta_0 + \beta_1 x$$

Now, we predict the odds of success

$$\log \left(\frac{P(x)}{1-P(x)} \right) = \beta_0 + \beta_1 x$$

exponentiating both sides

$$e^{\ln} \left(\frac{P(x)}{1-P(x)} \right) = e^{\beta_0 + \beta_1 x}$$

$$\frac{P(x)}{1-P(x)} = e^{\beta_0 + \beta_1 x}$$

$$\text{Let } y = e^{\beta_0 + \beta_1 x}$$

$$\text{then } \frac{P(x)}{1-P(x)} = y$$

$$P(x) = y / (1+y)$$

$$= y - y P(x)$$

$$P(x) + y P(x) = y$$

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

$$P(x) = \frac{y}{1+y}$$

$$P(x) = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$$

The equation of a sigmoid function

$$P(x) = \frac{e^{B_0 + B_2 x}}{1 + e^{B_0 + B_2 x}}$$

$$P(x) = \frac{1}{1 + e^{-(B_0 + B_1 x)}}$$

equation of Logistic Regression
that is Sigmoid function



Example of
sigmoid curve

Linear Regression

⇒ Used to solve Regression problems

⇒ The response variable are continuous in nature

⇒ It helps to calculate dependent variable where there is a change in independent variable

⇒ It is a straight line ⇒ An S-curve (S = Sigmoid)

⇒

Logistic Regression

⇒ Used to solve classification problem

⇒ The response variable is categorical in nature

⇒ It helps to calculate the possibility of a particular event taking place

⇒

Use-Case implementation of logistic Regression

Importing lib.

```
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
from sklearn import metrics
```

%matplotlib inline

```
digits = load_digits()
```

determining the total number of images & labels

```
print("Image Data shape", digits.data.shape)
print("Label Data shape", digits.target.shape)
```

Displaying some of the images & labels

```
plt.figure(figsize=(20, 4))
```

```
for index, (image, label) in enumerate(zip(digits.data[0:5], digits.target[0:5])):
    plt.subplot(1, 5, index + 1)
    plt.imshow(np.reshape(image, (8, 8)), cmap=plt.cm.gray)
    plt.title('Training : -i\n' + str(label), fontsize=20)
```

```
plt.show()
```

```
cmap=plt.cm.gray
```

```
plt.title('Training : -i\n' + str(label),  
         fontsize=20)
```

dividing dataset into training & test set

$x_train, x_test, y_train, y_test =$
`train_test_split(digit.data, digit.target,
test_size=0.23,
random_state=2)`

import the logistic Regression model

`from sklearn.linear_model import LogisticRegression`

Making an instance of the model &
training it

`logisticRegr = LogisticRegression()`

`logisticRegr.fit(x_train, y_train)`

Predicting output of 1st element of test set

`print(logisticRegr.predict(x_test[0].
reshape(1, -1)))`

Predicting the output of first 10 element of test

`logisticRegr.predict(x_test[0:10])`

Prediction for the entire dataset

`predictions = logisticRegr.predict(x_test)`

determining the accuracy of the model

score = logistic Regr. score (x-test, y-test)

print(score)

cm = metrics.confusion_matrix(y-test, prediction)

confusion matrix

plt.figure(figsize=(9,9))

sns. heat map (cm, annot=True, fmt=".3f",
line width=.5, square=True,
cmap='Blues_r');

plt.y_label('Actual label');

plt.x_label('Predicted label');

all-sample-title = 'Accuracy score : {}',
format(score)

plt.title(all-sample-title, size=15);

* A.N.O.V.A :- Analysis of Variance by R.A. Fisher

- ⇒ A statistical test for detecting difference in the group means when there is one parametric dependent variable & one or more independent variables.
- ⇒ It is used when we want to study 3 or more population study.

Example:-	NY	Texas	Oregon	
	18	18	21	$H_0: \mu_1 = \mu_2 = \mu_3$
	19	20	22	$H_a: \text{at least one mean is different}$
	20	16	17	
	21	20	18	
	22	21	22	
	23	20	19	
	18	18	21	
	19	19	20	
	20	17	18	
	21	13	23	
	\bar{x}_1	\bar{x}_2	\bar{x}_3	

Grand mean :- The mean of the sample means

$$\bar{\bar{x}} = \sum_{i=1}^k \left[\frac{n_i}{\sum_{j=1}^{n_i} x_{ij}} \right] \quad K = \text{population set } \{3\}$$

$$\sum n_i$$

n_i = sample size of i^{th} - population

$$\bar{\bar{x}} = \frac{584}{30} = 19.46667$$

Null hypothesis : A general statement that states that there is no relationship between two measured phenomena or no association among groups

Alternate hypothesis :- It states whenever something is happening, a new theory is preferred instead of an old one.

P value :- Probability of finding the observed, or more extreme, results when the null hypothesis of a study question is true. {P3}

Alpha value :- Criterion for determining whether a test statistic is statistically significant {α3}

LN	x^2	MN	x^2	LN	x^2
10	100	8	64	4	16
9	81	7	16	3	9
6	36	6	36	6	36
7	49	7	49	4	16
$\sum x_1$	$\frac{266}{266} =$	$\sum x_2 =$	$\frac{\sum x_2^2 = 165}{165}$	$\sum x_3 =$	$\frac{\sum x_3^2 = 77}{77}$
= 32	$\sum x_1^2$	25		17	

$$\text{Correction Term : } C_x \frac{(\sum x)^2}{n} = \frac{(32+25+17)^2}{12} = 456.33$$

$$\begin{aligned} \text{Sum of square of total : } SS_T &= \sum x^2 - C_x = \\ &= (266 + 165 + 77) - 456.33 \\ &= 51.67 \end{aligned}$$

Sum of square among groups

$$SS_A = \frac{(\sum X)^2}{n} - C_x = \left(\frac{32^2}{4} + \frac{25^2}{4} + \frac{17^2}{4} \right) - 465.33 \\ = 28.17$$

Sum of squares within groups

$$SS_W = SS_T = SSA \\ = 51.67 - 28.17 = 23.5$$

Mean of sum of squares among groups

$$\text{MSS}_A = \frac{SS_A}{k-1} = \frac{28.17}{3-1} = 14.085$$

Degree of freedom
among group

Mean of sum of squares within groups

$$\text{MSS}_W = \frac{SS_W}{N-k} = \frac{23.5}{12-3} = 2.611$$

Degree of freedom
within group

N ... Total no. of samples
k ... total categories or group

$$F \text{ Ratio} = \frac{\text{MSS}_A}{\text{MSS}_W} = \frac{14.085}{2.611} = 5.394$$

$$F(2, 9) < 5.294 \quad \alpha = 0.05 \quad \text{--- (1)} \\ F(2, 9) > 5.34 \quad \alpha = 0.01$$

Standard deviation σ

\Rightarrow Our vary from average. Low SD. means maximum data is near average & high S.D. means its not so

$$\text{Population} \quad \text{Sample}$$

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} \quad \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

$$\text{variance} = \frac{\sum (x - \bar{x})^2}{n}$$

A
Bessel's
Correction

① Reject null hypothesis & accept Alternative hypothesis

95% confidence

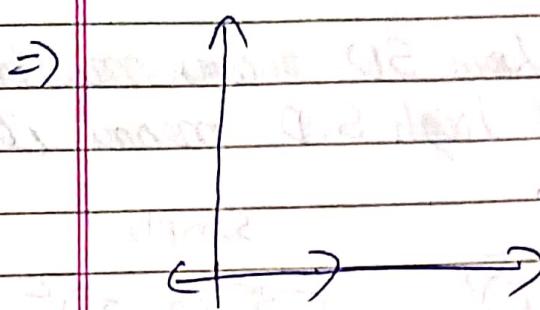
② Reject alternative hypothesis & accept null hypothesis

99% confidence

No : No significant effect of rain on number of questions solved.

H₁ : Significant effect of rain on Number of questions solved.

Cost Function [Continuation of Linear regression]



Idea is to choose θ_0, θ_1 so that $h_\theta(x)$ is close to y for our training example (x_i, y_i)

$$\Rightarrow \text{minimize } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=0}^m (h_\theta(x_i) - y_i)^2$$

$J(\theta_0, \theta_1)$

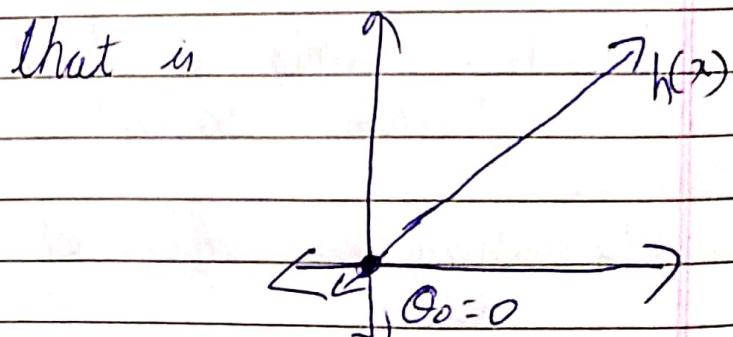
Cost function also squared error function

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

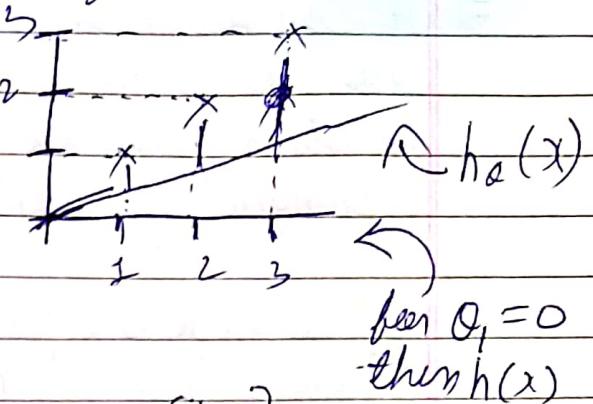
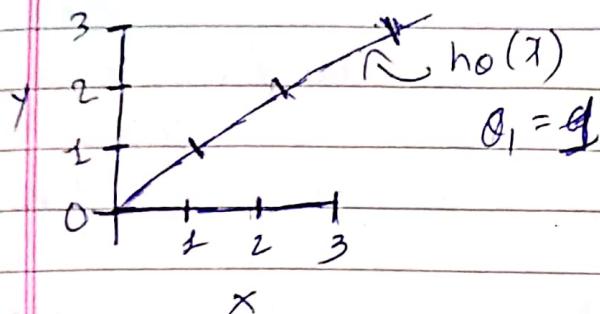
Parameter θ_0, θ_1

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

$$\Rightarrow \text{Example let consider } h_\theta(x) = \theta_1 x, \theta_0 = 0$$

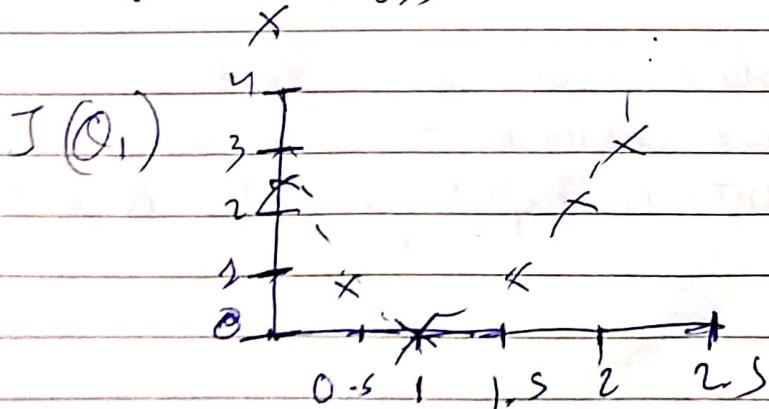


$h_0(x)$ for fixed θ , this is a function of x



$$\begin{aligned}
 J(\theta_0) &= \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 \\
 &= \frac{1}{2m} \sum_{i=1}^n (\theta_0 x^{(i)} - y^{(i)})^2 \\
 &= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0^2
 \end{aligned}$$

$J(\theta)$ function of the parameter θ ,



$$\begin{aligned}
 J(\theta) &= \frac{1}{2m} (1^2 + 2^2 + 3^2) \\
 &= 2 \cdot 3
 \end{aligned}$$

$$\begin{aligned}
 J(\theta_0) &= \frac{1}{2m} \left\{ (\theta_0 - 1)^2 + (2 - 2)^2 + (1 - 3)^2 \right\} \\
 &= \frac{1}{2 \times 3} (3.5) = 0.58
 \end{aligned}$$

- ⇒ A contour plot is a graph that contains many contour lines. A contour line of a two variable function has a constant value at all points of the same line. An example of such a graph is the one to the right below.
- ⇒ If we have one parameter hence can be described using 2D graph that is parameter vs $J(\theta)$ the cost function.
- ⇒ For more parameter dimension shall increase and hence difficult to visualize.

* Gradient descent

- ⇒ As a general function, for minimizing cost function in linear regression, also to minimize other function
- ⇒ We have $J(\theta_0, \theta_1)$ we want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Outline ⇒ Start with some θ_0, θ_1 ,

⇒ keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum

Algorithm

repeat until convergence

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

for $j = 0$ to $j = 1$

Correct Simultaneous update

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

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

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

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

α ... learning rate

* Multivariate Linear Regression

\Rightarrow Previously :- $h_\theta(x) = \theta_0 + \theta_1 x$

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

n = no. of features

$x^{(i)}$ = input features of i^{th} training example

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

* Clustering

- ⇒ Grouping data points & creating partitions based on similarity
- ⇒ Don't have response class

Types :- Centroid Based clustering
i) K-means

:- Connectivity Based Clustering

:- Distribution Based Clustering
:- Density Based Clustering

⇒ Sensitive to outliers

⇒ We need to supply the number of Cluster

Example :- Applications of Clustering

- Market Segmentation
- Social Network analysis
- Astronomical data analysis

K-Means

\Rightarrow Input :- No. of cluster that is value of K
 :- Training set

$$K = \{2, 3, 4, 10, 11, 12, 20, 25, 30\} \text{ for } K=2$$

$$m_1 = 4$$

$$m_2 = 12$$

$$K_1 = \{2, 3, 4\}$$

Note $m_1 = 3$

$$K_2 = \{10, 11, 12, 20, 25, 30\}$$

$$m_2 = 18 = \frac{108}{6}$$

$$K_1 = \{2, 3, 4, 10\}$$

$$m_1 = 4.75$$

$$K_2 = \{11, 12, 20, 25, 30\}$$

$$m_2 = 19.6$$

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

$$m_1 = 7$$

$$m_2 = 25$$

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$m_1 = 7$$

$$K_2 = \{20, 25, 30\}$$

$$m_2 = 25$$

\Rightarrow Thus we are getting same means.

K-mediods

i	x	y	$\textcircled{1} C_1 (3, y)$	$\textcircled{2} C_2 (7, y)$
x_1	2	6		
x_2	3	7		$ a-c + b-d $
x_3	3	8		
x_4	4	7		
x_5	6	2		
x_6	6	4		
x_7	7	3		
x_8	7	4		
x_9	8	5		
x_{10}	7	6		

$$\textcircled{1} \quad |2-3| + |6-4| \Rightarrow \textcircled{3}$$

$$0 + 4 \Rightarrow \textcircled{9}$$

$$1 + 3 \Rightarrow \textcircled{7}$$

$$3 + 2 \Rightarrow \textcircled{5}$$

$$3 + 0 \Rightarrow \textcircled{3}$$

$$4 + 1 \Rightarrow \textcircled{5}$$

$$5 + 1 \Rightarrow \textcircled{6}$$

$$4 + 2 \Rightarrow \textcircled{6}$$

$$\textcircled{2} \quad |2-7| + |6-4| \Rightarrow 7$$

$$7 + 4 \Rightarrow 8$$

$$3 + 3 \Rightarrow 6$$

$$1 + 2 \Rightarrow \textcircled{3}$$

$$1 + 0 \Rightarrow \textcircled{1}$$

$$0 + 1 \Rightarrow \textcircled{1}$$

$$1 + 1 \Rightarrow \textcircled{2}$$

$$0 + 2 \Rightarrow \textcircled{1}$$

Step II then cluster are

$$C_1: \{(2, 6), (3, 2), (4, 7), (3, 4)\}$$

$$C_2: \{(7, 4), (6, 2), (6, 4), (7, 3), (8, 5), (7, 6)\}$$

$$\text{Total Cost} = \sum_{i=1}^d |x_i - c_i|$$

$$= \{ \text{Cost } (3, 4), (2, 6) \quad \text{Cost } (7, 4) (8, 5) \\ \text{Cost } (3, 4) (3, 2) \quad \text{Cost } (7, 4) (7, 6) \\ \text{Cost } (3, 4) (4, 7) \quad ! \}$$

$$= (3+4+4) + (3+1+1+2+2)$$

$$= 20$$

Select one of non-medoids O'

Let $= (7, 3)$ i.e. x_7

∴ Non-medoids are $C(3, 4)$ & $O'(7, 3)$

	x	y	$(7, 3)$	$(3, 4)$	distances	dust	Total Cost
x_1	2	6	8	5	3	5	$\Rightarrow 3+5+5+$
x_3	3	8	9	5	4	5	$2+2+3+3$
x_4	4	7	7	5	3	4	$\Rightarrow 22$
x_5	6	2	5	3			
x_6	6	4	5	3			
x_7	7	4	4	3			
x_9	8	5	3	6			
x_{10}	7	6	3	6			