

12/01/2020

3.) Data Analysis in DEP

Date _____
Page No. _____

Distance Based Models K-means

Distance metric uses distance function which provides a relationship metric between each element of the data set.

There are many ways to reach a particular destination.

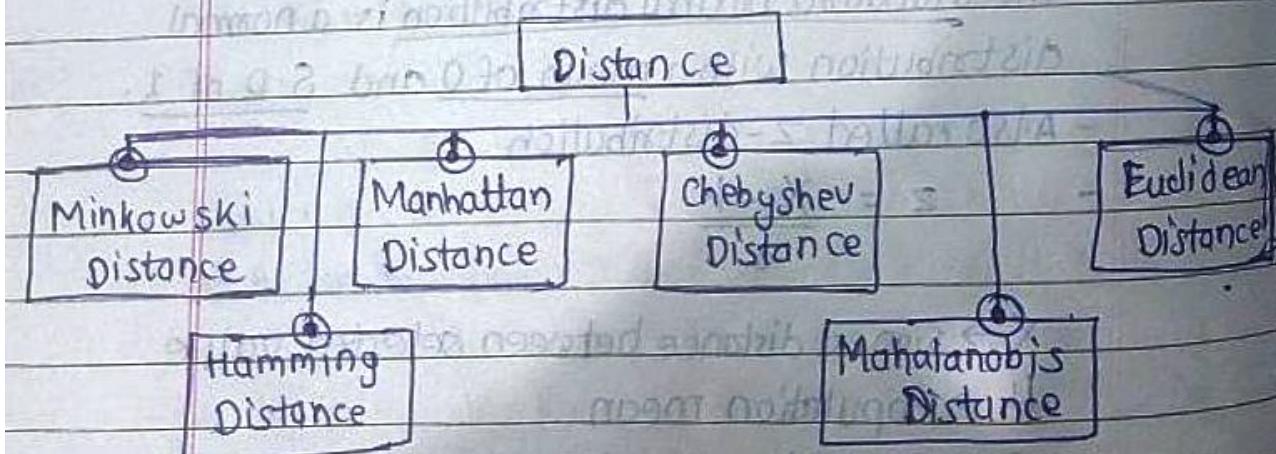
Distances are very useful to find similarities and dissimilarities between two instances.

Distance Measures w.r.t chess:

There are six different pieces to use while playing chess game. Each piece is governed to move by certain set of instructions which restricts its movement in some way. Likewise:

A mountain range can be considered as loose detour while using car, train or foot but it is easy to cross it flying.

Distances are used to measure similarity.



Minkowski Distance:

$$d(x, y) = \|x - y\|_m = \left(\sum_{i=1}^n (x_i - y_i)^m \right)^{1/m}$$

↓ Norm

If $m=1$, then it becomes Manhattan Distance
 $m=2$, Euclidean Distance
 $m=\infty$, Chebyshet Distance

Manhattan Distance:

Manhattan Distance is calculated as the sum of the absolute differences between two vectors.

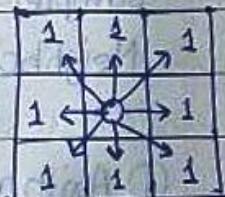
$$|x_1 - x_2| + |y_1 - y_2| = \|(x, y)\|_1$$

Euclidean Distance:

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Chebyshet Distance

$$Dis \propto (x, y) = \max_j |x_j - y_j|$$



Hamming Distance:

- Bits needed to change that is the Hamming distance
- Calculates the distance between two binary vectors, also referred to as binary strings or bits.

$$red = [1, 0, 0]$$

$$green = [0, 1, 0]$$

$$blue = [0, 0, 1]$$

Mahalanobis Distance:

$$d(x, y) = \|x - y\|_S^{-1} \|x - y\|$$

- Introduced by Prof. P. C. Mahalanobis in 1936 and is used in many statistical applications ever since.
- It transforms the columns into uncorrelated variables.
- Scale the columns into uncorrelated to make their variance equal to 1.
- Finally calculates Euclidean distance.

Properties of Distance:

- ① Distance between point and itself is zero.
- ② $\text{Dist}(x, y) \geq 0$.
- ③ $\text{Dist}(x, y) = \text{Dist}(y, x)$ are symmetric
- ④ Detours cannot shorten distance.

Exemplar:

- It is example, sample, instance.
- It can be ideal case for consideration for testing other sample.
- Neighbours :- closest of to exemplar.

① Arithmetic mean :

$$\frac{2+4}{2} = 3$$

② Geometric mean :

A term between two terms of geographic sequence is the geometric mean of two terms.

$$\text{Geometric mean of } 2 \text{ & } 18 = \sqrt{2 \times 18} = 6$$

③ Centroid:

- Centroid is an arithmetic mean

$$\mu = [1, 2, 3, 4, 5]$$

$$= 1+2+3+4+5/5$$

$$= 3$$

④ Medoid:

- Medoids are representative objects ~~one~~ of a data set or a cluster with a data set whose average dissimilarity to all objects in cluster is ~~minimum~~ minimal.

$$\mu = [1, 2, 3, 4, 5]$$

$$\text{Medoid} = 3$$

- It is always a member of given set.

Clustering:

- Grouping Unlabelled examples is called as clustering.

- If examples are labelled, then it becomes classification.

Distance Based Clustering:

Given a data Matrix X , the scatter matrix is,

$$S = (X - \mu)^T (X - \mu) = \sum_{i=1}^n (x_i - \mu)^T (x_i - \mu)$$

μ = row vector containing all columns means of X

Where scatter of X is defined as,

$$\text{Scatt}(X) = \sum_{i=1}^n \| (x_i - \mu)^2 \|$$

Imagine now that we have partitioned D into K subsets $D_1 \dots D_K = D$, and let μ_j denote the mean of D_j . Let S be the scatter matrix of D , and S_j be the scatter matrices of D_j . These scatter matrices have following relationship:

$$S = \sum_{j=1}^K (S_j + B)$$

Here, B is the scatter matrix that results by replacing each point in D with the corresponding μ_j .

K-means Clustering:

The K-means problem is NP complete, which means that there is no efficient solution to find the global minimum and we need to resort to the heuristic algorithm.

The algorithm iterates between partitioning the data using the nearest-centroid decision rule, and calculating centroids from a partition.

Algorithm:

Input: Data $d \subseteq \mathbb{R}^d$; no. of clusters $K \in \mathbb{N}$

Output: K cluster means $\mu_1, \dots, \mu_K \in \mathbb{R}^d$.

① Randomly initialize K vectors $\mu_1, \dots, \mu_K \in \mathbb{R}^d$.

② Repeat

③ Assign each $x \in D$ to $\arg \min_j \text{Dis}_2(x, \mu_j)$

④ for $j = 1$ to K do

⑤ $D_j \leftarrow \{x \in D \mid x \text{ assigned to cluster } j\}$;

⑥ $\mu_j = \frac{1}{|D_j|} \sum_{x \in D_j} x$

⑦ end

⑧ until no change in μ_1, \dots, μ_K ,

⑨ return μ_1, \dots, μ_K

Ex: Given Data Points are

$(2,9), (1,5), (7,4), (5,8), (6,4), (7,5), (1,2), (4,9)$

Let,

$$x_1 = (2,9) \quad x_2 = (1,5) \quad x_3 = (7,4) \quad x_4 = (5,8)$$

$$x_5 = (6,4) \quad x_6 = (7,5) \quad x_7 = (1,2) \quad x_8 = (4,9)$$

Let us consider $K=3$ — (any number) —

Initial centroids, $c_1 = x_1(2,9)$ $c_2 = x_4(5,8)$

$$c_3 = x_7(1,2)$$

— (can take any 3 at random) —

Iteration 1

	Distance from c_1 $x_1(2,9)$	Distance from $c_2(5,8)$	Distance from $c_3(1,2)$
$x_1(2,9)$	$\sqrt{(2-2)^2 + (9-9)^2} = 0$	$\sqrt{(2-5)^2 + (9-8)^2} = 3.16$	$\sqrt{(2-1)^2 + (9-2)^2} = 7.07$
$x_2(1,5)$	$\sqrt{(1-2)^2 + (5-9)^2} = 4.12$	$\sqrt{(1-5)^2 + (5-8)^2} = 5$	$\sqrt{(1-1)^2 + (5-2)^2} = 3$
$x_3(7,4)$	$\sqrt{(7-2)^2 + (4-9)^2} = 7.07$	$\sqrt{(7-5)^2 + (4-8)^2} = 4.47$	$\sqrt{(7-1)^2 + (4-2)^2} = 6.32$
$x_4(5,8)$	$\sqrt{(5-2)^2 + (8-9)^2} = 3.16$	$\sqrt{(5-5)^2 + (8-8)^2} = 0$	$\sqrt{(5-1)^2 + (8-2)^2} = 7.21$
$x_5(6,4)$	$\sqrt{(6-2)^2 + (4-9)^2} = 6.40$	$\sqrt{(6-5)^2 + (4-8)^2} = 4.12$	$\sqrt{(6-1)^2 + (4-2)^2} = 5.38$
$x_6(7,5)$	$\sqrt{(7-2)^2 + (5-9)^2} = 6.40$	$\sqrt{(7-5)^2 + (5-8)^2} = 3.60$	$\sqrt{(7-1)^2 + (5-2)^2} = 6.70$
$x_7(1,2)$	$\sqrt{(1-2)^2 + (2-9)^2} = 7.07$	$\sqrt{(1-5)^2 + (2-8)^2} = 7.21$	$\sqrt{(1-1)^2 + (2-2)^2} = 0$
$x_8(4,9)$	$\sqrt{(4-2)^2 + (9-9)^2} = 2$	$\sqrt{(4-5)^2 + (9-8)^2} = 1.41$	$\sqrt{(4-1)^2 + (9-2)^2} = 7.61$

* Closest ones are highlighted.

New Centroids

$$c_1 = (3,9)$$

$$c_2 = \cancel{(5,8)} \rightarrow (6.25, 5.25)$$

$$c_3 = (1,3.5)$$

Iteration 2

Page No. 50
Date

	Distance from $C_1(3, 9)$	Distance from $C_2(6.25, 5.25)$	Distance from $C_3(1, 3.5)$
$x_1(2, 9)$	$\sqrt{(2-3)^2 + (9-9)^2} = 1$	$\sqrt{(2-6.25)^2 + (9-5.25)^2} = 5.73$	$\sqrt{(2-1)^2 + (9-3.5)^2} = 5.5$
$x_2(1, 5)$	$\sqrt{(1-3)^2 + (5-9)^2} = 4.4$	$\sqrt{(1-6.25)^2 + (5-5.25)^2} = 5.81$	$\sqrt{(1-1)^2 + (5-3.5)^2} = 1.5$
$x_3(7, 4)$	$\sqrt{(7-3)^2 + (4-9)^2} = 6.40$	$\sqrt{(7-6.25)^2 + (4-5.25)^2} = 1.45$	$\sqrt{(7-1)^2 + (4-3.5)^2} = 6.0$
$x_4(5, 8)$	$\sqrt{(5-3)^2 + (8-9)^2} = 2.23$	$\sqrt{(5-6.25)^2 + (8-5.25)^2} = 3.02$	$\sqrt{(5-1)^2 + (8-3.5)^2} = 6.0$
$x_5(6, 4)$	$\sqrt{(6-3)^2 + (4-9)^2} = 5.83$	$\sqrt{(6-6.25)^2 + (4-5.25)^2} = 1.27$	$\sqrt{(6-1)^2 + (4-3.5)^2} = 5.0$
$x_6(7, 5)$	$\sqrt{(7-3)^2 + (5-9)^2} = 5.65$	$\sqrt{(7-6.25)^2 + (5-5.25)^2} = 0.75$	$\sqrt{(7-1)^2 + (5-3.5)^2} = 6.0$
$x_7(1, 2)$	$\sqrt{(1-3)^2 + (2-9)^2} = 7.28$	$\sqrt{(1-6.25)^2 + (2-5.25)^2} = 6.17$	$\sqrt{(1-1)^2 + (2-3.5)^2} = 1.5$
$x_8(4, 9)$	$\sqrt{(4-3)^2 + (9-9)^2} = 1$	$\sqrt{(4-6.25)^2 + (9-5.25)^2} = 4.37$	$\sqrt{(4-1)^2 + (9-3.5)^2} = 6.28$

New Centroids:

$$C_1 = (3, 9) \quad C_2 = (6.25, 5.25) \quad C_3 = (1, 3.5)$$

Iteration 3 is same as Iteration 2

Therefore above centroids are required output.

Advantages:

- It is very easy

- It is done in polynomial time

Disadvantages:

- sensitive to outliers, an object with extremely large value may substantially distort the distribution of data.

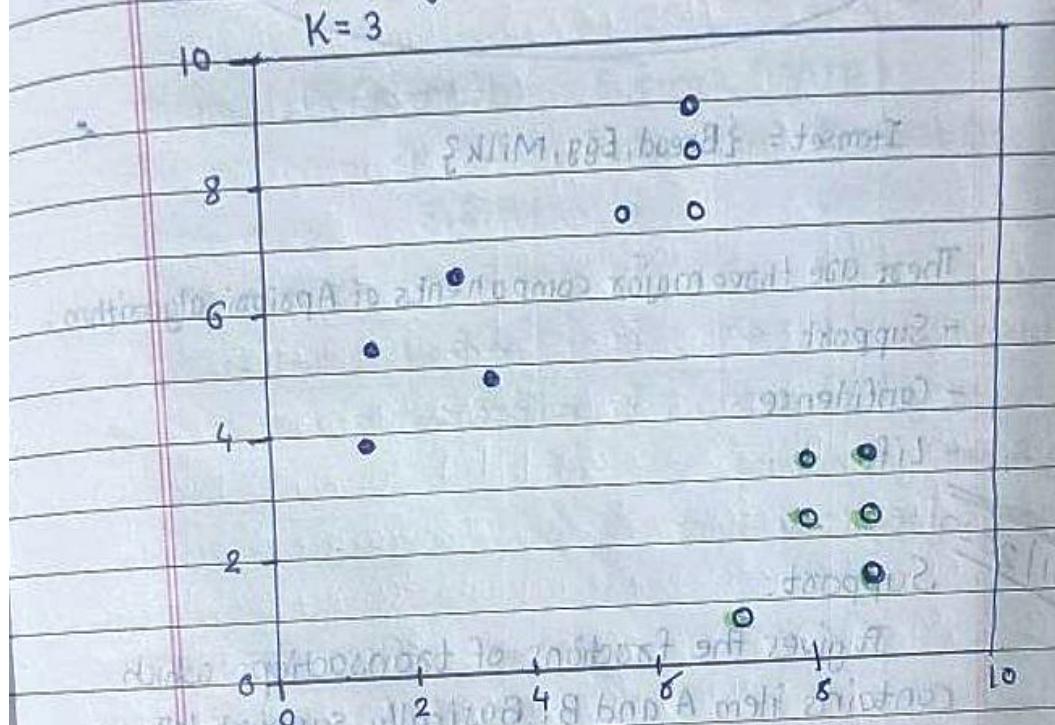
- K is not known in advance

- Cannot handle categorical data.

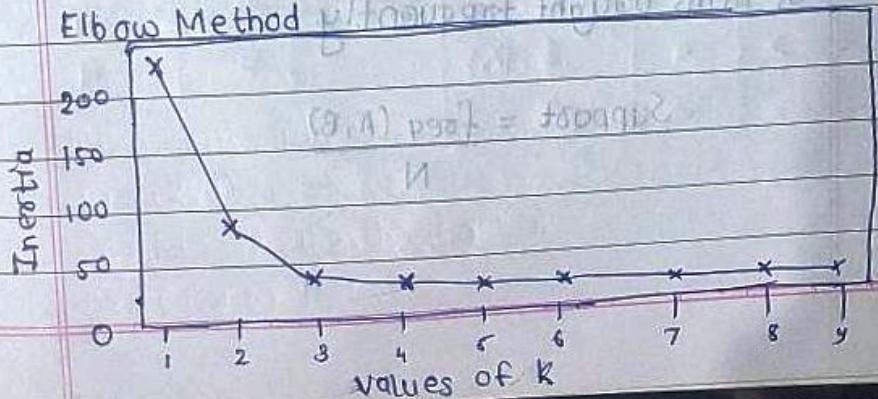
Determining K in K-means:

- Deciding K in K-means is difficult as data is large and unlabelled.
- Distortion:- It is calculated as an average of the squared distances from the cluster centers of the respective clusters. Typically, the Euclidean distance metric is used.
- Inertia :- It is the sum of squared distances of samples to their closest cluster center.

Scatter plot of given data set



Elbow Method



Association rule:

- Association rules analysis is a technique to uncover how items are associated to each other.
- The process of identifying associations between products is called as association rule mining.

Apriori :-

$\{ \text{Bread, Eggs} \} \rightarrow \{ \text{Milk} \}$

Antecedent Consequent

Itemset = {Bread, Egg, Milk}

These are three major components of Apriori algorithm:

- Support
- Confidence
- Lift

01/2020

Support:

It gives the fractions of transactions which contains item A and B. Basically Support tells us about the frequently bought items or the combination of items bought frequently.

Support = $\frac{\text{freq}(A, B)}{N}$

Confidence:

It tells us how often the items A and B occurs together, given the number of times A occurs.

$$\text{Confidence} = \frac{\text{freq}(A, B)}{\text{freq}(A)}$$

Lift:

Lift indicates the strength of a rule over the random occurrence of A and B. It is basically tells us the strength of any rule.

$$\text{Lift}(A \rightarrow B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A) \cdot \text{Support}(B)}$$

Support

↳ Freq. of an item in
given set of transaction

Confidence

↳ It is conditional
probability of
purchasing item B given
A is already purchased.

Lift

↳ It is same as confidence

but also considers the

popularity of item B also

$$\text{Lift}(A \rightarrow B) = \text{Conf}(A \text{ and } B) / \text{Supp}(B)$$

Approach:-

- Set a minimum value for support and confidence.
- Extract all the subsets having higher value of support than minimum threshold.
- Select all rules from the subsets with confidence value higher than minimum threshold.
- Order the rules by descending order of Lift.

Association Rule Mining:-

Ex:-

Given:- Support = 2 Confidence = 60%

Sr. No.	Tid	Items
1.	T ₁	1, 3, 4
2	T ₂	2, 3, 5
3	T ₃	1, 2, 3, 5
4	T ₄	2, 5
5	T ₅	1, 3, 5

→ Frequency Count

Item	Frequency	Item	Frequency
1	3	1	3
2	3	2	3
3	4	3	4
As Support=2 →	4	5	4
	5		
	4		

- Two Set Items

Item Set	Frequency		Item Set	Frequency
1, 2	1 (X)		1, 3	3
1, 3	3		1, 5	2
1, 5	2	→	2, 3	2
2, 3	2		2, 5	3
2, 5	3		3, 5	3
3, 5	3			

- Three Set Items

Item Set	Frequency		Item Set	Frequency
1, 3, 5	2	→	1, 3, 5	2
1, 2, 3	1		2, 3, 5	2
2, 3, 5	2			

- Setting Rules:

① 1, 3, 5

↪ 1^3 → 5

$$\text{Confidence} = \text{Support}(1, 3, 5) / \text{Support}(1, 3)$$

$$= 2/3$$

$$(2/3) / (2/3) = 66.66\%$$

↪ 3^5 → 1

$$\text{Confidence} = \text{Support}(1, 3, 5) / \text{Support}(1, 3)$$

$$= 2/3$$

$$= 66.66\%$$

↪ 1 → 3^5

$$\text{Confidence} = \text{Support}(1, 3, 5) / \text{Support}(1)$$

$$= 2/3$$

$$= 66.66\%$$

$\hookrightarrow 3 \rightarrow 1^5$

$$\begin{aligned}\text{Confidence} &= \text{Support}(1,3,5) / \text{Support}(3) \\ &= 2/4 \\ &= 50\% \rightarrow (\text{Rejected})\end{aligned}$$

$\hookrightarrow 5 \rightarrow 1^3$

$$\begin{aligned}\text{Confidence} &= \text{Support}(1,3,5) / \text{Support}(5) \\ &= 2/4 \\ &= 50\% \rightarrow (\text{Rejected})\end{aligned}$$

② $2,3,5$

$\hookrightarrow 2^3 \rightarrow 5$

$$\begin{aligned}\text{Confidence} &= \text{Support}(2,3,5) / \text{Support}(2,3) \\ &= 2/2 \\ &= 100\%\end{aligned}$$

$\hookrightarrow 2^5 \rightarrow 3$

$$\begin{aligned}\text{Confidence} &= \text{Support}(2,3,5) / \text{Support}(2,5) \\ &= 2/3 \\ &= 66.66\%\end{aligned}$$

$\hookrightarrow 3^5 \rightarrow 2$

$$\begin{aligned}\text{Confidence} &= \text{Support}(2,3,5) / \text{Support}(3,5) \\ &= 2/3 \\ &= 66.66\%\end{aligned}$$

$\hookrightarrow 2 \rightarrow 3^5$

$$\begin{aligned}\text{Confidence} &= \text{Support}(2,3,5) / \text{Support}(2) \\ &= 2/3 \\ &= 66.66\%\end{aligned}$$

↳ $3 \rightarrow 2^3 5$

$$\begin{aligned}\text{Confidence} &= \text{Support}(2,3,5) / \text{Support}(3) \\ &= 2/4 \\ &= 50\% \text{ (Rejected)}\end{aligned}$$

↳ $5 \rightarrow 2^5 3$

$$\begin{aligned}\text{Confidence} &= \text{Support}(2,3,5) / \text{Support}(5) \\ &= 2/4 \\ &= 50\% \text{ (Rejected)}\end{aligned}$$

18/01/2021

Regression:

$(x) f = y$

Supervised
Learning

Classification Regression

Regression is used to predict real value or it can also be termed as Estimator function. Regression is used to predict Stock Price, Weather forecast and home prices in a particular area.

Regression analysis:-

A set of statistical processes for estimating the relationships between a dependent variable (outcome variable) and one or more independent variables (predictors). The most common form of regression analysis is linear regression.

- A function estimator, regressor, is mapping $f: X \rightarrow Y$.
The regression learning problem is to learn a function estimator from example $(x_i, f(x_i))$. For instance, we might want to learn an estimator for the Dow Jones index, or the FTSE 100 based on selected economic indicators.

(a) Fig 2.1 (e, s) - $\hat{Y} = b_0 + b_1 X$

- In Regression, our aim is to show relationship between two variable (independent variables) X and Response variable (target/dependent variable) Y .

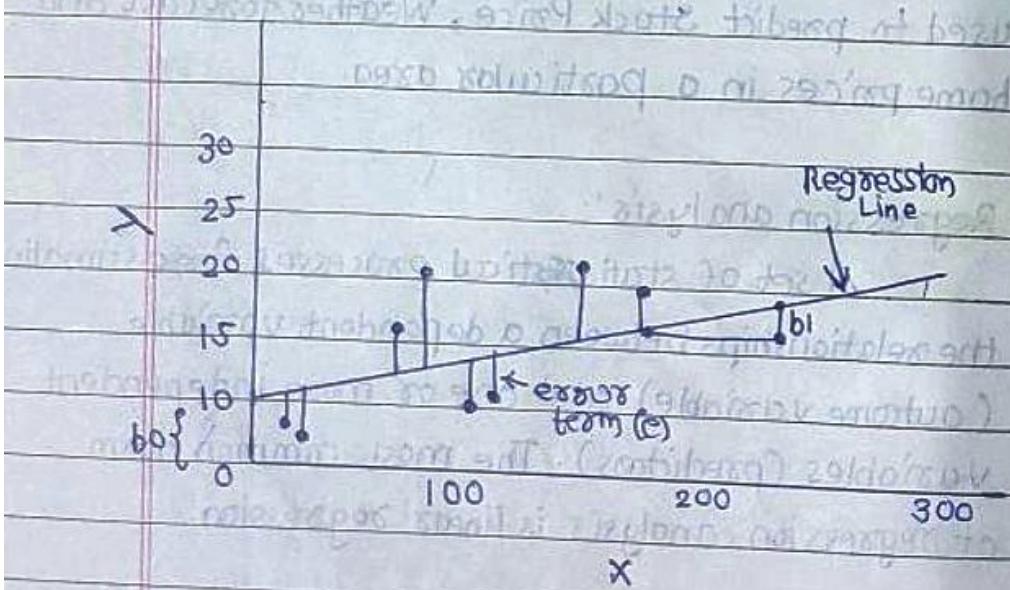
$$Y = f(X)$$

But while building a relation between dependent variable X and independent variable y for the function may / may not give exact value for y .

$$Y = f(X)$$

$$\hat{Y} = f(X)$$

$$\text{Residual error} = Y - \hat{Y}$$



$\hat{y} \rightarrow y$ Predicted

Page No.	59
Date	

①	②	③	④	⑤	⑥	⑦	⑧	⑨
$(x - \bar{x})^2$	x	y	X-mean	Y-mean	$(x - \bar{x})^2$	\hat{y}	$y - \hat{y}$	$(y - \hat{y})^2$
1.44	1	2	-2	-2	4	2.8	-0.8	0.64
0.36	2	4	-1	0	1	3.4	0.6	0.36
0	3	5	0	1	0	4	1	1
0.36	4	4	1	0	1	4.6	-0.6	0.36
1.44	5	5	2	1	4	5.2	-0.2	0.04
3.6	[3]	[4]			10		2.4	6
	↑	↑						
	x-	y-	mean	mean				

$$\text{Eqn of line : } Y = mx + c$$

$$Y = b_0 + b_1 x$$

b_0 - Intercept

b_1 = Slope

$$b_1 = \frac{\text{sum of all}(x - \text{mean}) \times (y - \text{mean})}{\text{sum of all}(x - \text{mean})^2}$$

$$b_0 = Y_{\text{mean}} - b_1 * X_{\text{mean}}$$

$$= \frac{(-2) \times (-2)}{10}$$

$$= (-2 \times -2) + (-1 \times 0) + (0 \times 1) + (1 \times 0) + (2 \times 1) / 10$$

$$= 4 + 0 + 0 + 2 / 10$$

$$= 6 / 10$$

$$\therefore = 0.6$$

$$b_0 = 4 - (0.6) \times 3 = 2.2$$

$$\hat{y} = 2.2 + 0.6x$$

$$SSE = \sum (y - \hat{y})^2 = 2.4$$

$$SSR = \sum (\hat{y} - \bar{y})^2 = 3.6$$

$$SST = SSR + SSE = 2.4 + 3.6 = 6$$

$$R^2 = SSR / SST = 3.6 / 6 = 0.6$$

Implementation Using python:

[1]

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
import pandas as pd
```

[2]

```
data = pd.DataFrame ({'x': [1,2,3,4,5], 'y': [2,4,5,4,5]})
```

[3]

```
print(data)
```

	x	y
0	1	2
1	2	4
2	3	5
3	4	4
4	5	5

[4]

```
x-mean = np.mean(data['x'])
```

```
y-mean = np.mean(data['y'])
```

```
n = len(data)
```

[5]

```
print(x-mean)
```

3.0

[6]

```
print(y-mean)
```

4.0

[7]

```
x = data['x'].values
```

```
y = data['y'].values
```

[8]

```
print(x)
```

[1,2,3,4,5]

[9] `print [y]
[2, 4, 5, 4, 5]`

[10] `num = 0
den = 0
for i in range (n):
 num += ((X[i] - X_mean) * (Y[i] - Y_mean))
 den += ((X[i] - X_mean)) ** 2
m = num / den
c = Y_mean - (m * X_mean)`

[11] `X = np.linspace (1, 10, 2)
Y = (m * X) + c
plt.plot (X, Y, color = 'blue', label = 'Regression Line')
plt.scatter (X, Y, color = 'green', label = 'Scatter')
plt.xlabel ('X')
plt.ylabel ('Y')
plt.legend ()
plt.show ()`

plots

[12] `SSo = 0
SST = 0
SSE = 0
for i in range (n):
 data ['Y-P'] = (m * data ['X']) + c
 Y_Pred = data ['Y-P']
 for i in range (n):
 SSE += (Y[i] - Y_Pred[i]) ** 2
 SSo += (Y_mean - Y_Pred) ** 2
SST = SSo + SSE
R_squared = SSo / SST
print (SSo, SSE, SST)
print (R_squared)`

[13] Point data

x	y	y-p̂
0	1	2
1	2	4
2	3	5
3	4	4.6
4	5	5.2

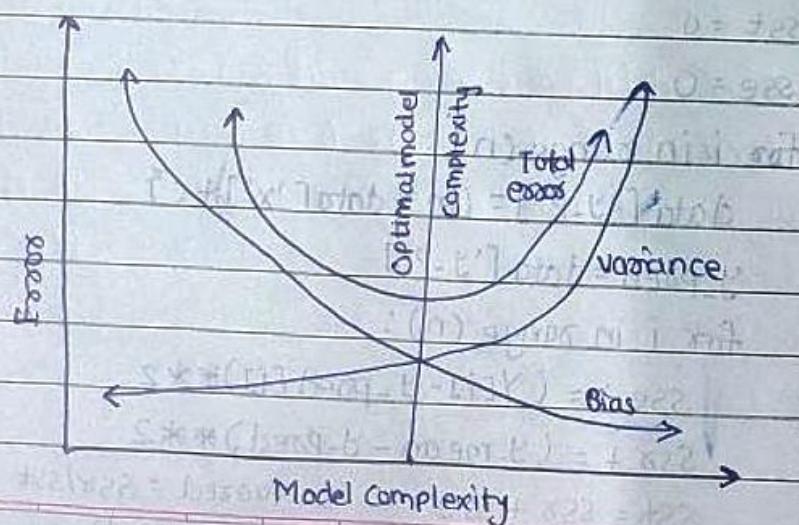
[14]

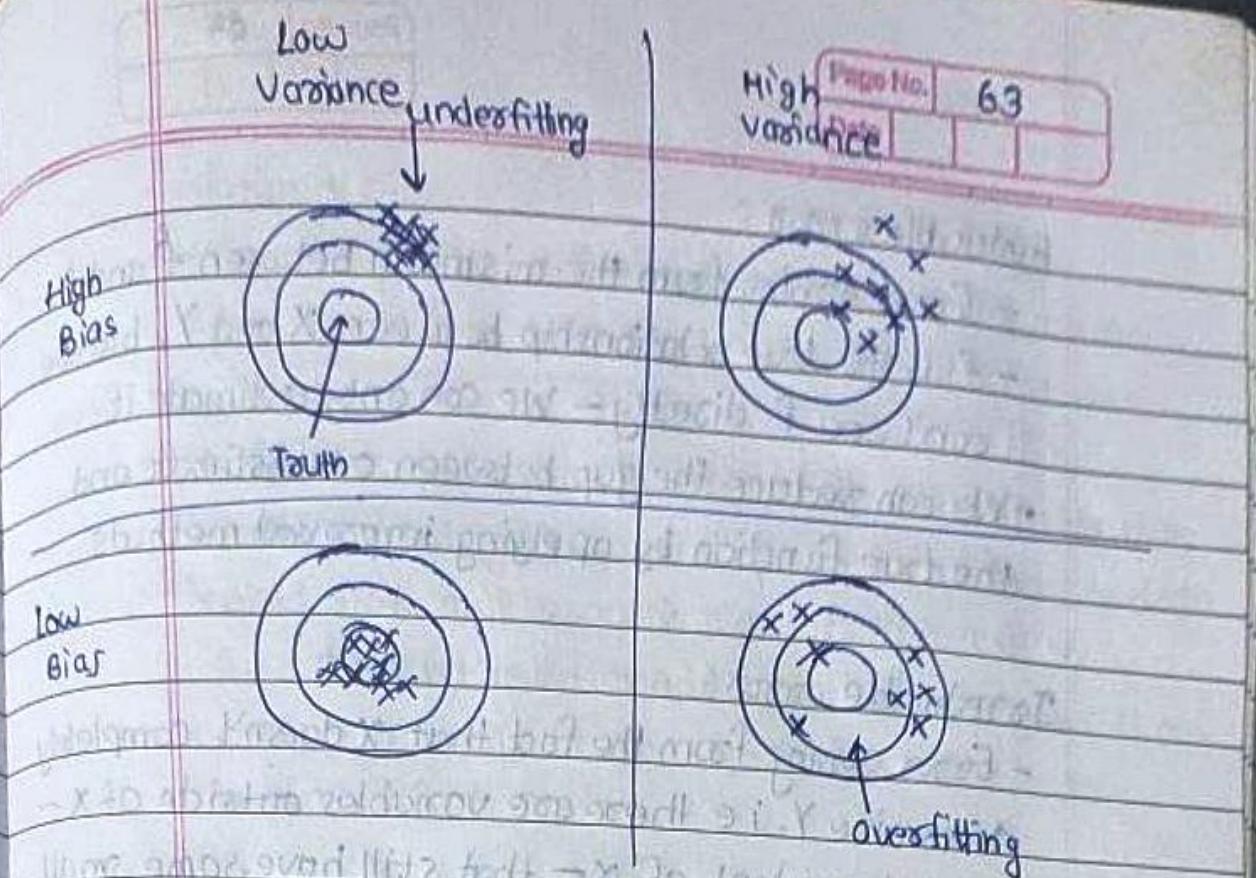
```

plt.scatter(X,y, color = 'green', label = 'scatter')
plt.plot([x,y], color = 'blue', label = "Regression Line")
plt.axhline(y = y_mean, color = 'orange', linestyle = '-')
plt.axvline(x = x_mean, color = 'orange', linestyle = '-')
plt.plot([X[i], X[i], Y[i], y_pred[i]], color = 'red', linestyle = '-')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.show()

```

Bias-Variance Tradeoff :





If our model is too simple and has very few parameters then it may have high bias and low variance. On the other hand, if our model has large numbers of parameters, then its going to have high variance and low bias.

So, we need to find the right / good balance without overfitting or underfitting data.

This tradeoff in complexity is where there is a tradeoff between bias and variance. An algorithm can't be more complex and less complex at the same time.

Errors

$$Y = f(x) + e$$

e = Reducible error + Irreducible error

Reducible errors :

- Errors arising from the mismatch between f and \hat{f} .
- f is the true relationship between X and Y , but we can't see f directly - We can only estimate it.
- We can reduce the gap between our estimate and the true function by applying improved methods.

Irreducible errors :

- Error arising from the fact that X doesn't completely determine Y . i.e there are variables outside of X - and independent of X - that still have some small effect on Y .
- The only way to improve prediction errors related to irreducible errors is to identify these outside influences and incorporate them as predictors.

$$\text{Reducible errors} = \text{Bias} + \text{Variance}$$

Bias :

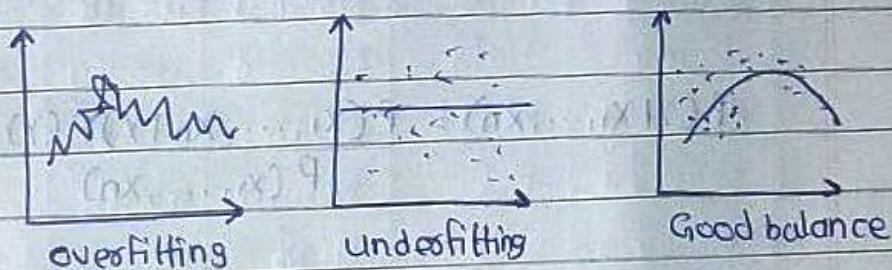
- The difference between the average prediction of our model and the correct value which we are trying to predict.
- Model with high bias pays very little attention to the training data and oversimplifies the model.
- It always leads to high errors on training and test data.

Variance:

- Variability of model prediction for a given data point or a value which tells us spread of our data.
- Model with high variance pays a lot of attention to the training data and does not generalize data which it hasn't seen before.
- As a result, such models perform very well on training data but has high error rates on test data.

Error in prediction:

$$Error(x) = E[(Y - \hat{f}(x))^2]$$



Underfitting:

- happens when a model is unable to capture the underlying pattern of the data.
- These models usually have ~~high~~ variance and high bias.
- It happens when we have very less amount of data to build an accurate model.

Overfitting :

- happens when our model captures the noise along with the underlying pattern of data.
- It happens when we train our model a lot over noisy dataset. These models have low bias and high variance.
- These models are very complex like Decision trees which are prone to overfitting.

~~22/11/2020~~

Naive Bayes :

- An probabilistic Approach which will identify the posterior probability of occurrence of an event

The Bayes Classifier :

$$P(Y|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y)P(Y)}{P(X_1, \dots, X_n)}$$