

# Problem\_ 1: Clustering

## Table of Contents

### Contents

Executive Summary.....	2
Introduction .....	2
Data Description .....	2
Sample of the dataset.....	2
Exploratory Data Analysis.....	2
Let us check the types of variables in the data frame. ....	3
Check for missing values in the dataset.....	3
Q 1.1 Read the data, do the necessary initial steps, and exploratory data analysis (Univariate, Bi-variate, and multivariate analysis).....	4
Q 1.2 Do you think scaling is necessary for clustering in this case? Justify .....	7
Q1.3 Apply hierarchical clustering to scaled data. Identify the number of optimum clusters using Dendrogram and briefly describe them.....	8
Q1.4 Apply K-Means clustering on scaled data and determine optimum clusters. Apply elbow curve and silhouette score. Explain the results properly. Interpret and write inferences on the finalized clusters.....	10
Q1.5 Describe cluster profiles for the clusters defined. Recommend different promotional strategies for different clusters.....	13

## Executive Summary

A leading bank wants to develop a customer segmentation to give promotional offers to its customers. They collected a sample that summarizes the activities of users during the past few months. You are given the task to identify the segments based on credit card usage.

## Introduction

The purpose of this whole exercise is to explore the dataset. Do the exploratory data analysis. Explore the dataset and analysis of clustering techniques convert data are group wise customer past month credit card use .Past month credit card spending summarize data bank received payment and min limit payment amount , max spent in single shopping 1000s and min are 100s. Using hierarchical clustering and K-men for find group of customer performance.

## Data Description

1. Spending	: Amount spent by the customer per month (in 1000s)
2. advance payments	: Amount paid by the customer in advance by cash (in 100s)
3. probability_of_full_payment	: Probability of payment done in full by the customer to the bank
4. current balance	: Balance amount left in the account to make purchases (in 1000s)
5. credit limit	: Limit of the amount in credit card (10000s)
6. min_payment_amt	: minimum paid by the customer while making payments for purchases made monthly (in 100s)
7. max_spent_in_single_shopping	: Maximum amount spent in one purchase (in 1000s)

## Sample of the dataset

	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amt	max_spent_in_single_shopping
0	19.94	16.92	0.8752	6.675	3.763	3.252	6.550
1	15.99	14.89	0.9064	5.363	3.582	3.336	5.144
2	18.95	16.42	0.8829	6.248	3.755	3.368	6.148
3	10.83	12.96	0.8099	5.278	2.641	5.182	5.185
4	17.99	15.86	0.8992	5.890	3.694	2.068	5.837

Fig\_0.1

**Conclusion | insight:** This is Sample of given data all information is shown in this manner

## Exploratory Data Analysis

```
RangeIndex: 210 entries, 0 to 209
Data columns (total 7 columns):
 #   Column            Non-Null Count  Dtype  
 --- 
 0   spending          210 non-null    float64
 1   advance_payments  210 non-null    float64
 2   probability_of_full_payment  210 non-null    float64
 3   current_balance   210 non-null    float64
 4   credit_limit      210 non-null    float64
 5   min_payment_amt  210 non-null    float64
 6   max_spent_in_single_shopping 210 non-null    float64
dtypes: float64(7)
memory usage: 11.6 KB
```

**Fig\_0.2**

**Conclusion | insight:** Seven type of information in this data seat spending, advance payments probability\_of\_full\_payment ,current\_balance ,credit limit and min\_payment\_amt ,max\_spent\_in\_single\_shopping.data is 11.6 kb and 0 to 210 row and 7 columns.

## Check for missing values in the dataset

```
spending          0
advance_payments 0
probability_of_full_payment 0
current_balance 0
credit_limit     0
min_payment_amt 0
max_spent_in_single_shopping 0
dtype: int64
```

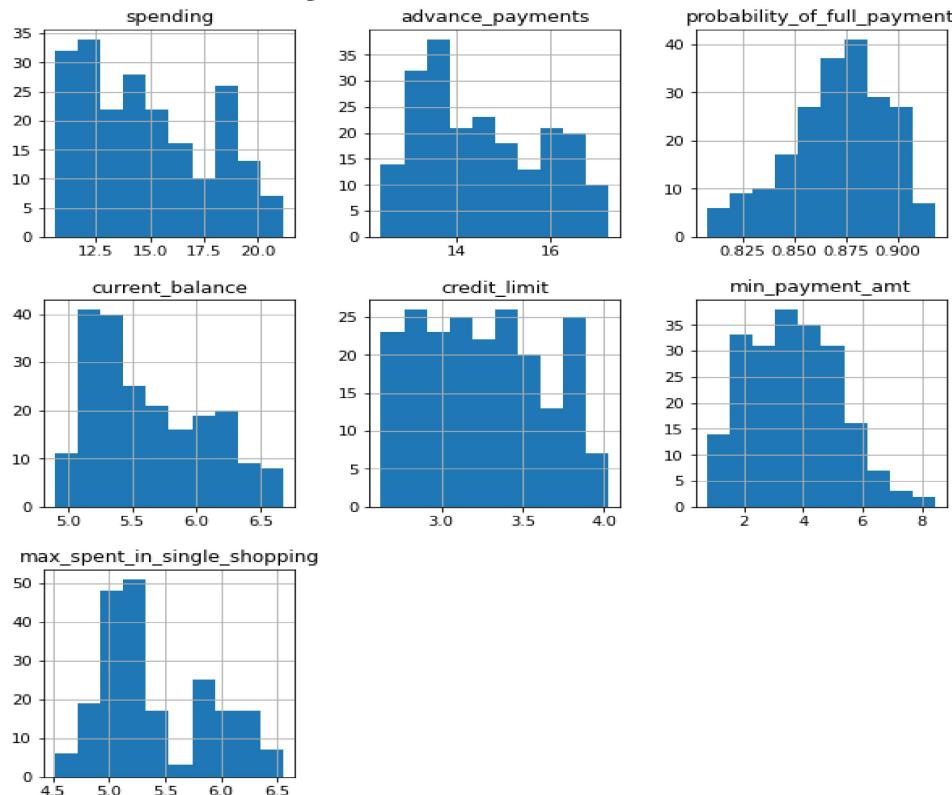
**Fig\_0.3**

**Conclusion | insight:** From the above results we can see that there is no missing value present in the dataset.

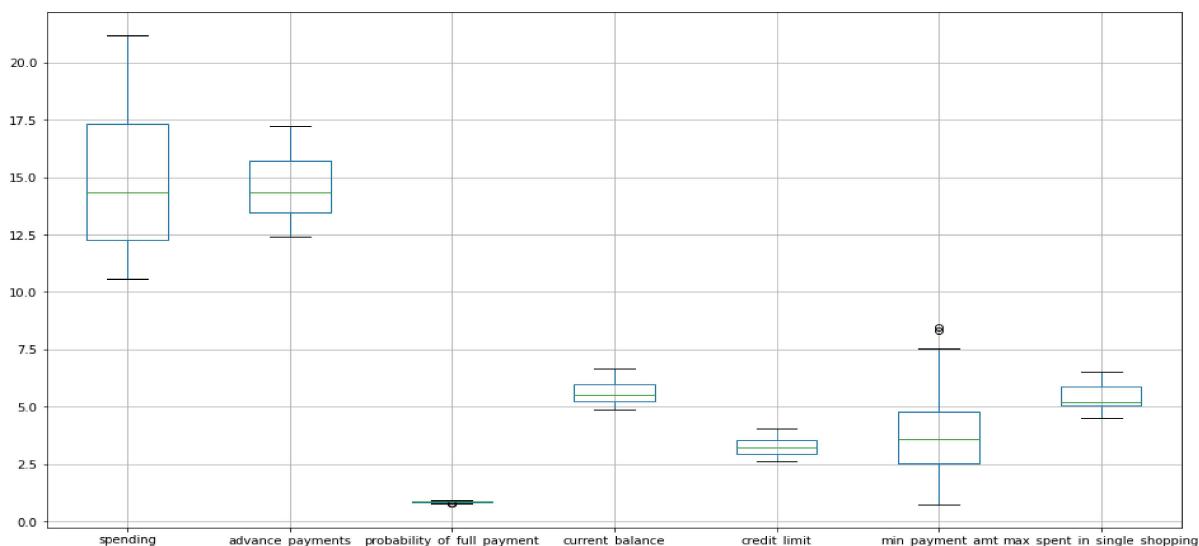
Q 1.1 Read the data, do the necessary initial steps, and exploratory data analysis (Univariate, Bi-variate, and multivariate analysis).

## Solution-

### Univariate analysis



Fig\_0.4

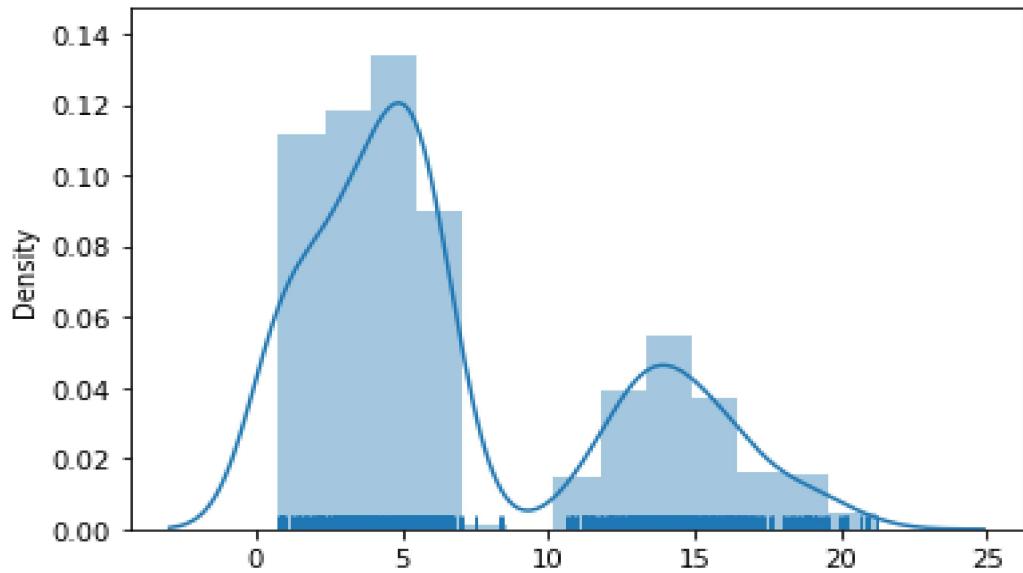


Fig\_0.5

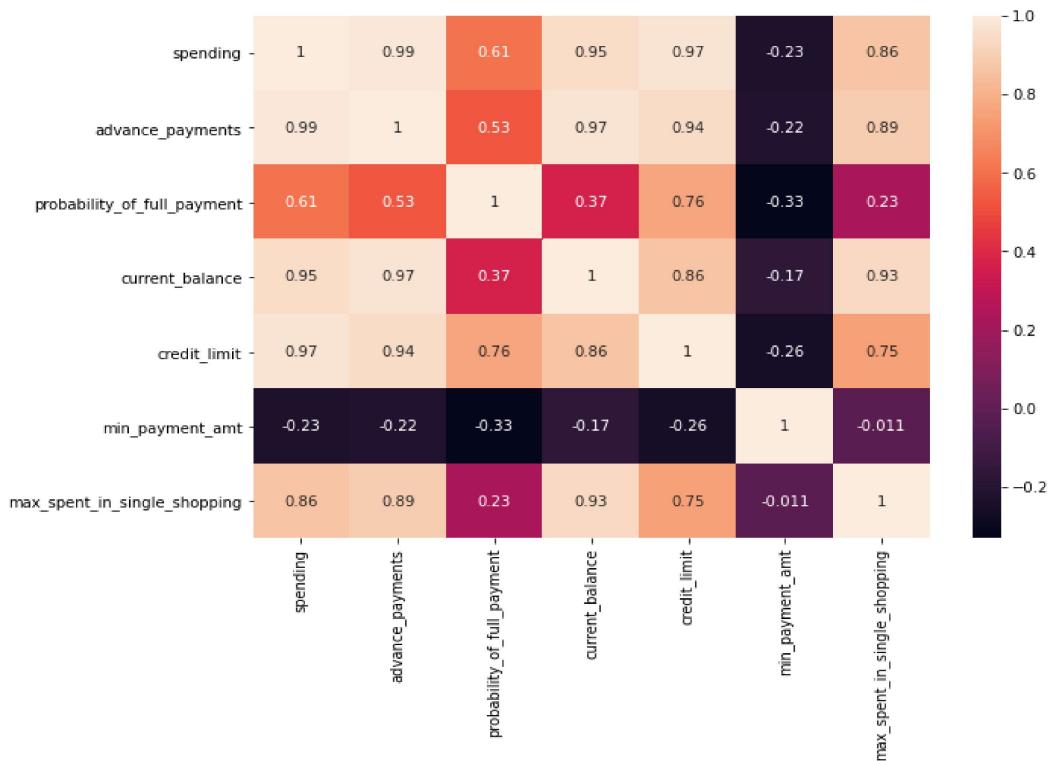
**Conclusion | insight:** All columns data single behaviour show in Fig\_0.4 & Fig\_0.5.

## Bivariate Analysis

Column- spending graph



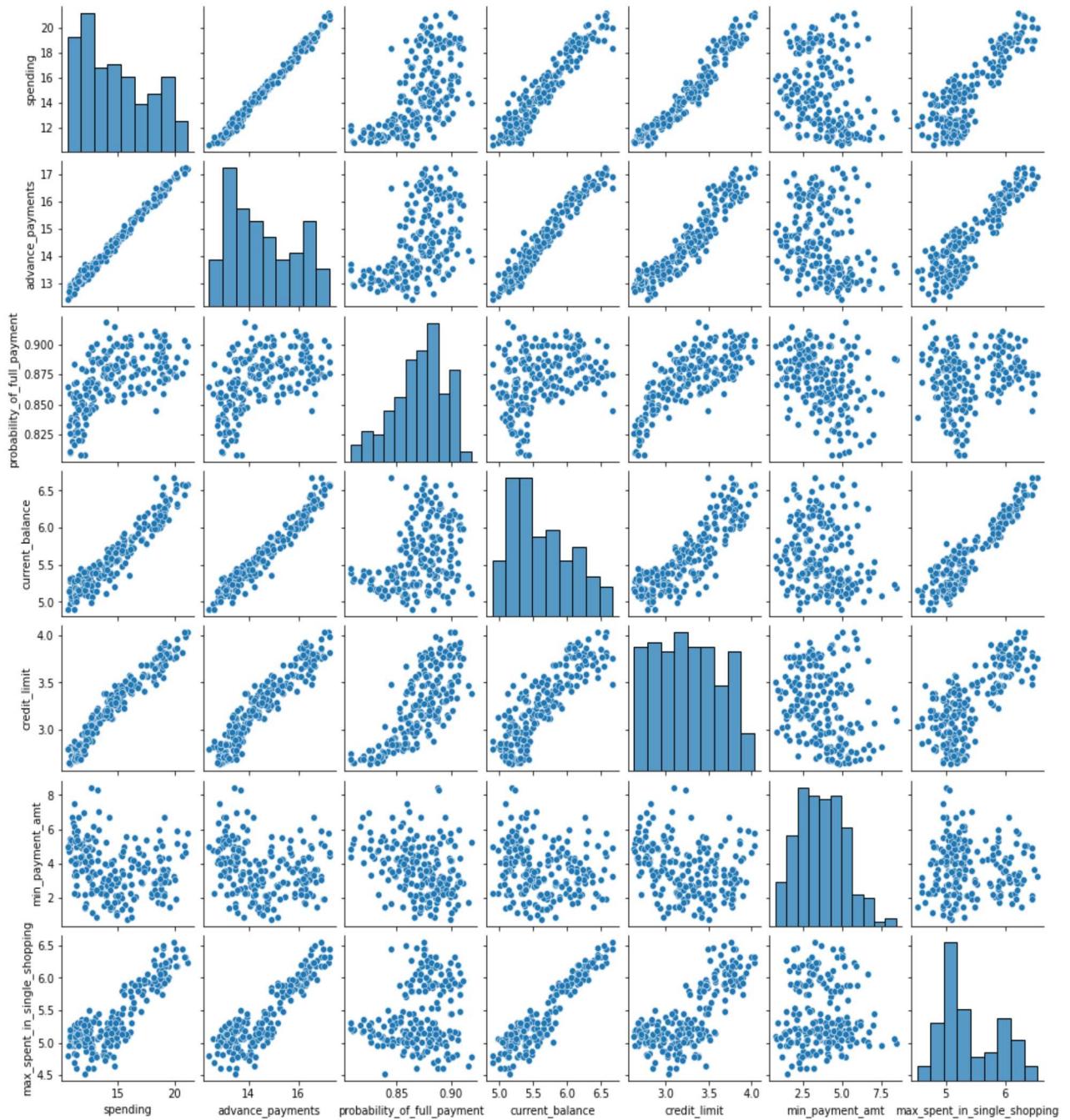
**Fig\_0.6**



**Fig\_0.7**

**Conclusion | insight:** Correlation of all two data behaviour in show in fig\_0.7.

## Multivariate Analysis



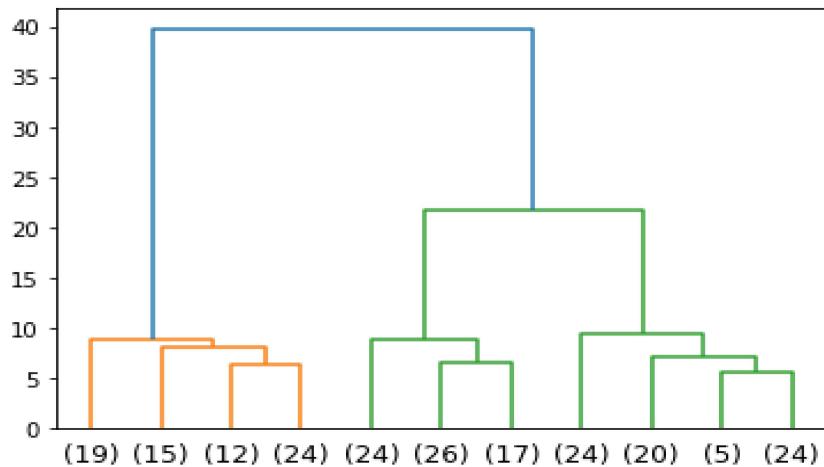
**Fig\_0.8**

**Conclusion | insight:** All data pair correction and behaviour show in fig\_0.8.





## Dendrogram of ward method using 10 pointer graph



**Fig\_11**

I have Used Ward linkage method to link the features of the dataset. This is truncated Dendrogram with P= 10. Since the original dendrogram was not clear at x-axis, so it was advisable to truncate it and get a neater diagram. We can also calculate the number of observations in each cluster by summing the numbers at x-axis. From dendrogram it is clear that no of optimum clusters should be 2.

## Hierarchical clustering value count of clusters

```

1    70
2    67
3    73
Name: clusters, dtype: int64

```

**Fig\_12**

Doing Cluster profiling for the 3 clusters as obtained from Hierarchical clustering clustering and taking their mean and generating the data frame for profiles.

## Final result add clusters in main data

	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amt	max_spent_in_single_shopping	clusters
0	19.94	16.92	0.8752	6.675	3.763	3.252	6.550	1
1	15.99	14.89	0.9064	5.363	3.582	3.336	5.144	3
2	18.95	16.42	0.8829	6.248	3.755	3.368	6.148	1
3	10.83	12.96	0.8099	5.278	2.641	5.182	5.185	2
4	17.99	15.86	0.8992	5.890	3.694	2.068	5.837	1

**Fig\_13**

Cluster 1: Customer with high spendings, and having less probability of full payment.

Cluster 2: Customer with Medium spendings and having high probability of full payment

Cluster 3: Customer with low spendings and having high probability of full payment

**Conclusion | insight:** Based on the dataset, 3 clusters can be optimum to do customer segmentation as it is also viable from the dendrogram in Full Payer and Revolvers. Because, Non payers are not available in the dataset.

**Q1.4** Apply K-Means clustering on scaled data and determine optimum clusters. Apply elbow curve and silhouette score. Explain the results properly. Interpret and write inferences on the finalized clusters.

## Solution-

### KMeans clusters three labels

```
[2, 0, 2, 1, 2, 1, 1, 0, 2, 1, 2, 0, 1, 2, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1,
2, 1, 0, 2, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 2, 2, 0, 2, 2,
1, 1, 0, 2, 2, 2, 1, 2, 2, 2, 2, 2, 1, 1, 1, 2, 0, 1, 1, 0, 0, 2, 2,
2, 0, 2, 1, 0, 1, 2, 2, 1, 2, 0, 1, 2, 0, 0, 0, 0, 2, 1, 0, 2, 0,
2, 1, 0, 2, 0, 1, 1, 2, 2, 2, 1, 2, 0, 2, 0, 2, 0, 2, 2, 1, 1, 2,
0, 0, 2, 1, 1, 2, 0, 0, 1, 2, 0, 1, 1, 1, 0, 0, 2, 1, 0, 0, 1, 0,
0, 2, 1, 2, 2, 1, 2, 0, 0, 0, 1, 1, 0, 1, 2, 1, 0, 1, 0, 1, 0, 1, 0, 0,
1, 0, 0, 1, 0, 2, 2, 1, 2, 2, 2, 1, 0, 0, 0, 1, 0, 1, 0, 2, 2, 2,
0, 1, 0, 1, 0, 0, 0, 0, 2, 2, 1, 0, 0, 1, 1, 0, 1, 2, 0, 2, 2, 1,
2, 1, 0, 2, 0, 1, 2, 0, 2, 0, 0, 0])
```

**Fig\_12**

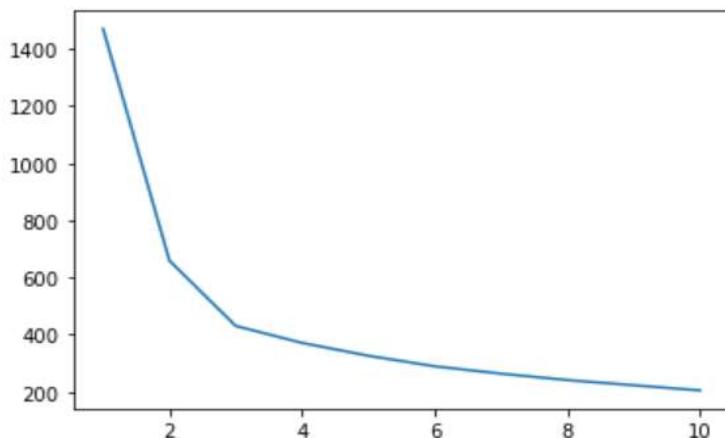
Magnitude of all the data is on same scale. There is not a large difference between the magnitude of higher and lower values which essentially is essence of scaling.

## Using K - Elbow Method

	0
<b>0</b>	1470.000000
<b>1</b>	659.171754
<b>2</b>	430.658973
<b>3</b>	371.301721
<b>4</b>	326.367602
<b>5</b>	289.420125
<b>6</b>	263.699963
<b>7</b>	241.319111
<b>8</b>	223.121257
<b>9</b>	205.184198

Fig\_13

## WSS Range graph as per K - Elbow Method



Fig\_14

Let's go for Silhouette score to verify the optimum no of clusters, if the silhouette score for  $k=2$  i.e. 2 cluster is better than 3 clusters ( $k=3$ ) then optimum number of clusters will be 2 otherwise 3.

## Silhouette score Cluster evaluation for KMeans 2 clusters

Silhouette score 2 clusters	0.46577247686580914
-----------------------------	---------------------

## Silhouette score Cluster evaluation for KMeans 3 clusters

Silhouette score 3 clusters	0.4007270552751299
-----------------------------	--------------------

## KMeans clustering value count of clusters

```
0    71
1    72
2    67
Name: Clus_kmeans, dtype: int64
```

Fig\_0.15

## Final result add clusters in main data

spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amt	max_spent_in_single_shopping	clusters	Clus_kmeans
19.94	16.92	0.8752	6.675	3.763	3.252	6.550	1	2
15.99	14.89	0.9064	5.363	3.582	3.336	5.144	3	0
18.95	16.42	0.8829	6.248	3.755	3.368	6.148	1	2
10.83	12.96	0.8099	5.278	2.641	5.182	5.185	2	1
17.99	15.86	0.8992	5.890	3.694	2.068	5.837	1	2

Fig\_0.16

Doing Cluster profiling for the 3 clusters as obtained from k-means clustering and taking their mean and generating the data frame for 3 profiles.

Cluster 0: Customer with high spendings, and having less probability of full payment.

Cluster 1: Customer with Medium spendings and having high probability of full payment

Cluster 2: Customer with low spendings and having high probability of full payment

**Conclusion | insight:** Based on the dataset, 3 clusters can be optimum to do customer segmentation as it is also viable from the k-mean in Full Payer and Revolvers. Because, Non payers are not available in the dataset. Elbow Method is using to find of behave of cluster.

**Q1.5** Describe cluster profiles for the clusters defined. Recommend different promotional strategies for different clusters.

## Solution-

### Recommendation for different promotional Strategies:

1. All result shown in both method Bank should provide more lucrative offers (like cashback /discounts on purchase on different websites/products) to low spending customers to convert them in high spending customers.
2. Almost data was as Bank can also promote some cashback on full payment mechanism which will motivate customers to do more full payments instead of minimum payments.
3. Bank should increase credit limits of full\_payments /advance\_payments customers to motivate them to do more spending and max spending in single shopping, which may result in EMI and bank can earn on it.
4. Bank can lower the EMI processing fees for high spending customers so that customer can go for EMI options and bank can earn interest on EMI.
5. Using credit Card almost costumer are timely pay one group cluster 1 & cluster 2 are hay ratio so cannulisation are best .

# Problem\_ 2: CART-RF-ANN

## Table of Contents

### Contents

Executive Summary.....	15
Introduction .....	15
Data Description .....	15
Sample of the dataset.....	15
Exploratory Data Analysis.....	16
Let us check the types of variables in the data frame. ....	3
Check for missing values in the dataset.....	3
Check for Duplicate values in the dataset.....	
Q 2.1 Read the data, do the necessary initial steps, and exploratory data analysis (Univariate, Bi-variate, and multivariate analysis).....	18
Q 2.2 Data Split: Split the data into test and train, build classification model CART, Random Forest, Artificial Neural Network.....	22
Q 2.3 Performance Metrics: Comment and Check the performance of Predictions on Train and Test sets using Accuracy, Confusion Matrix, Plot ROC curve and get ROC_AUC score, classification reports for each model.....	26
Q 2.4 Final Model: Compare all the models and write an inference which model is best/optimized.....	33
Q 2.5 Inference: Based on the whole Analysis, what are the business insights and recommendations.....	35

## Executive Summary

An Insurance firm providing tour insurance is facing higher claim frequency. The management decides to collect data from the past few years. You are assigned the task to make a model which predicts the claim status and provide recommendations to management. Use CART, RF & ANN and compare the models' performances in train and test sets.

## Introduction

The purpose of this whole exercise is to explore the dataset. Do the exploratory data analysis. Explore the dataset and analysis of Use CART, RF & ANN and compare the models' performances in train and test sets. The management decides to collect data from the past few years. Target of claimed data all modal are perform as per target.

## Data Description

1. Claimed	Claim Status - Target
2. Agency Code	Code of tour firm
3. Type	Type of tour insurance firms
4. Channel	Distribution channel of tour insurance agencies
5. Product	Name of the tour insurance products
6. Duration	Duration of the tour
7. Destination	Destination of the tour
8. Sales	Amount of sales of tour insurance policies
9. Commission	The commission received for tour insurance firm
10. Age	Age of insured

## Sample of the dataset

	Age	Agency_Code	Type	Claimed	Commision	Channel	Duration	Sales	Product Name	Destination
0	48	C2B	Airlines	No	0.70	Online	7	2.51	Customised Plan	ASIA
1	36	EPX	Travel Agency	No	0.00	Online	34	20.00	Customised Plan	ASIA
2	39	CWT	Travel Agency	No	5.94	Online	3	9.90	Customised Plan	Americas
3	36	EPX	Travel Agency	No	0.00	Online	4	26.00	Cancellation Plan	ASIA
4	33	JZI	Airlines	No	6.30	Online	53	18.00	Bronze Plan	ASIA

Fig\_1.0

**Conclusion | insight:** This is Sample of given data all information is shown in this manner.

## Exploratory Data Analysis

```
RangeIndex: 3000 entries, 0 to 2999
Data columns (total 10 columns):
 #   Column      Non-Null Count  Dtype  
 --- 
  0   Age         3000 non-null    int64  
  1   Agency_Code 3000 non-null    object  
  2   Type         3000 non-null    object  
  3   Claimed     3000 non-null    object  
  4   Commision    3000 non-null    float64 
  5   Channel      3000 non-null    object  
  6   Duration     3000 non-null    int64  
  7   Sales        3000 non-null    float64 
  8   Product Name 3000 non-null    object  
  9   Destination   3000 non-null    object  
 dtypes: float64(2), int64(2), object(6)
 memory usage: 234.5+ KB
```

Fig\_1.2

**Conclusion | insight:** Nine type of information in this data seat Age, agency\_code , Type , Claimed , commission, Channel, Duration ,sales, product Name and Destination data is 234.5 kb and 0 to 2999 row and 10 columns.

## Check for missing values in the dataset

```
Age          0
Agency_Code  0
Type         0
Claimed     0
Commision    0
Channel      0
Duration     0
Sales        0
Product Name 0
Destination   0
dtype: int64
```

Fig\_1.3

**Conclusion | insight:** From the above results we can see that there is no missing value present in the dataset.

## Check for Duplicate values in the dataset

```
ins_df.duplicated().value_counts()
```

```
False    2861  
True     139  
dtype: int64
```

```
ins_df.duplicated().sum()
```

```
139
```

```
ins_df.drop_duplicates(inplace=True)
```

```
ins_df.duplicated().value_counts()
```

```
False    2839  
dtype: int64
```

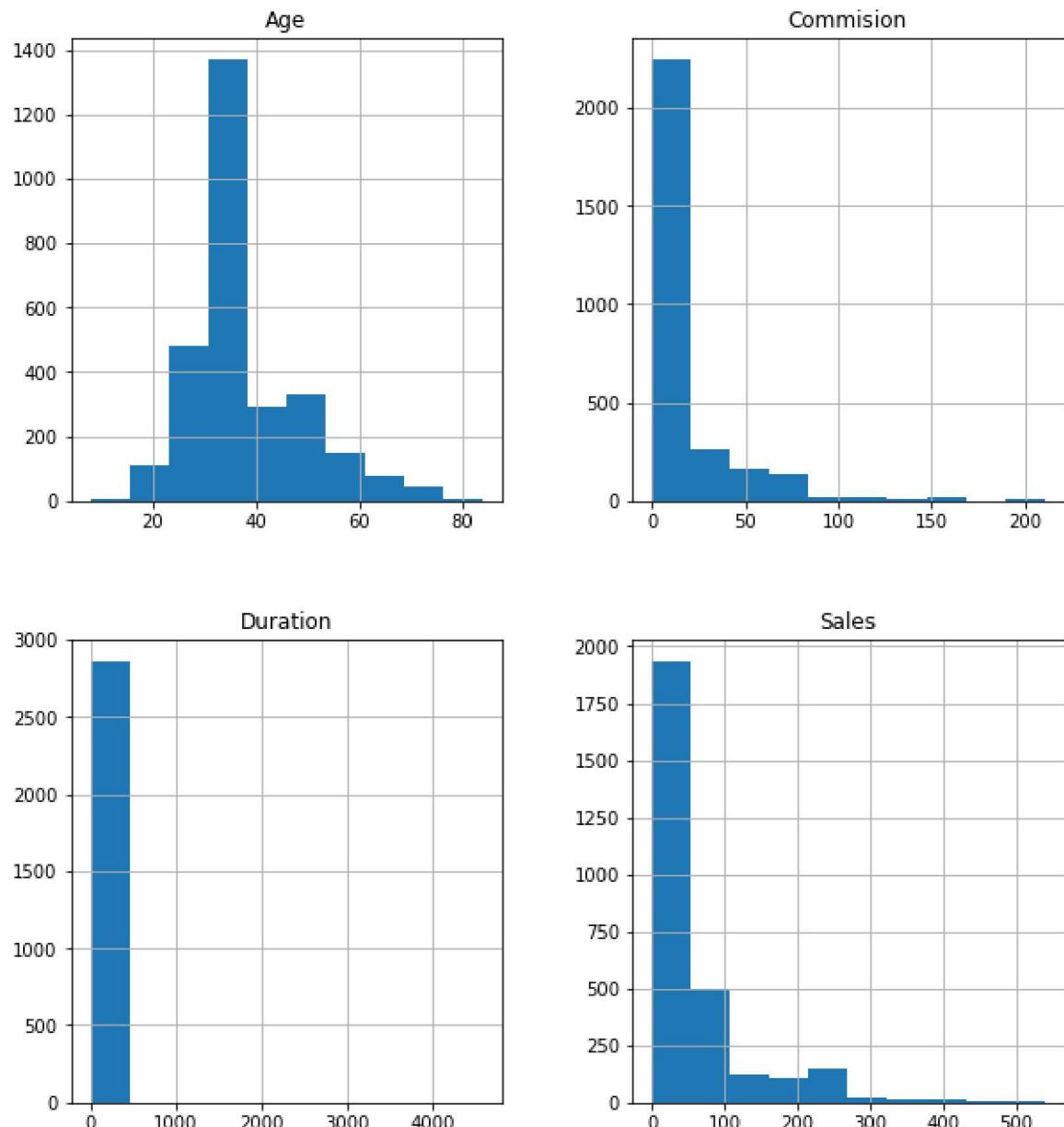
Fig\_1.4

**Conclusion | insight:** In this data 139 duplicate value identified but after trite no any duplicate value in this data.

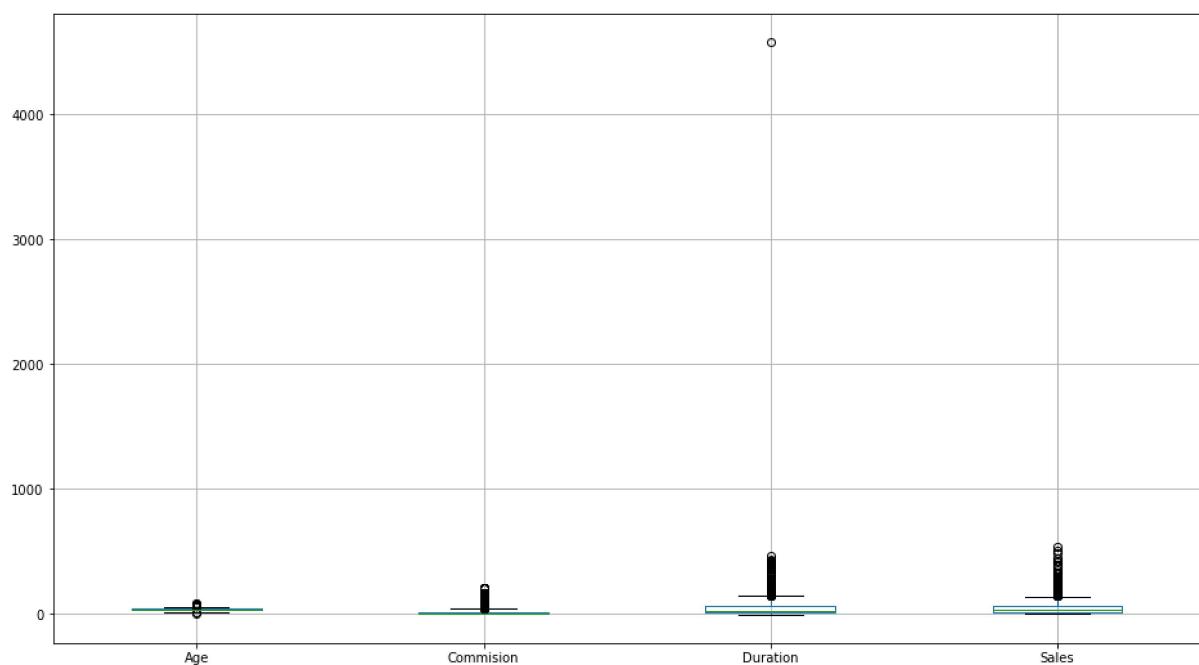
**Q 2.1** Read the data, do the necessary initial steps, and exploratory data analysis (Univariate, Bi-variate, and).

## Solution-

### Univariate analysis



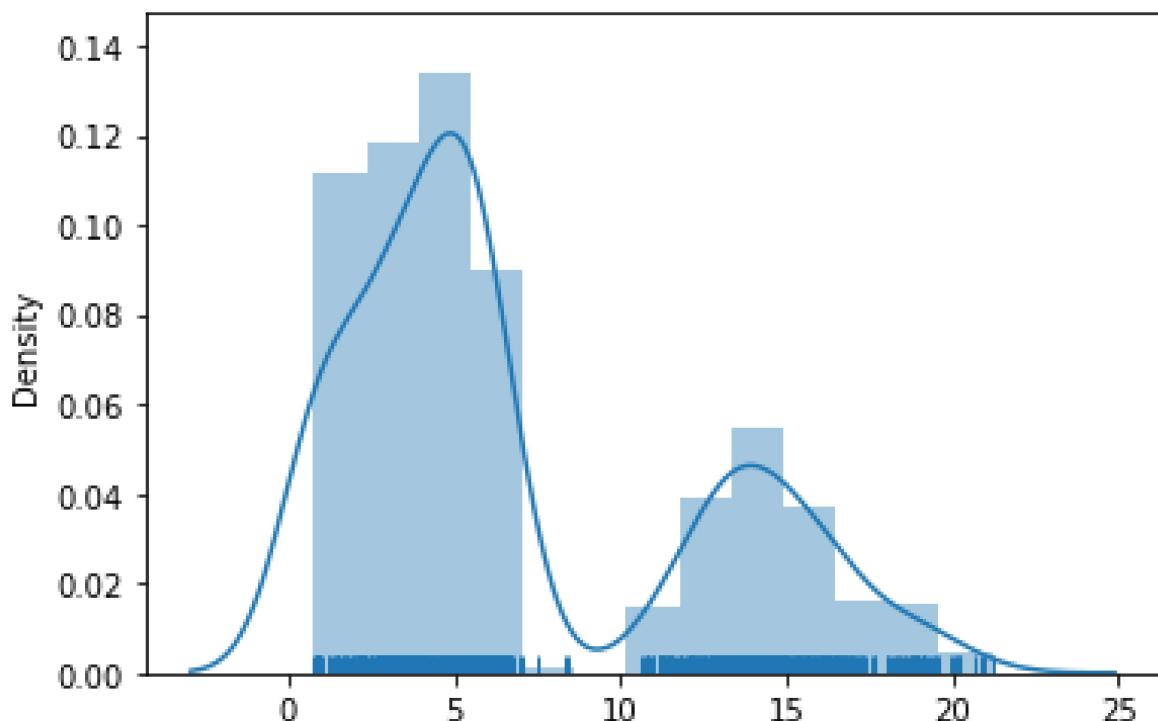
**Fig\_1.5**

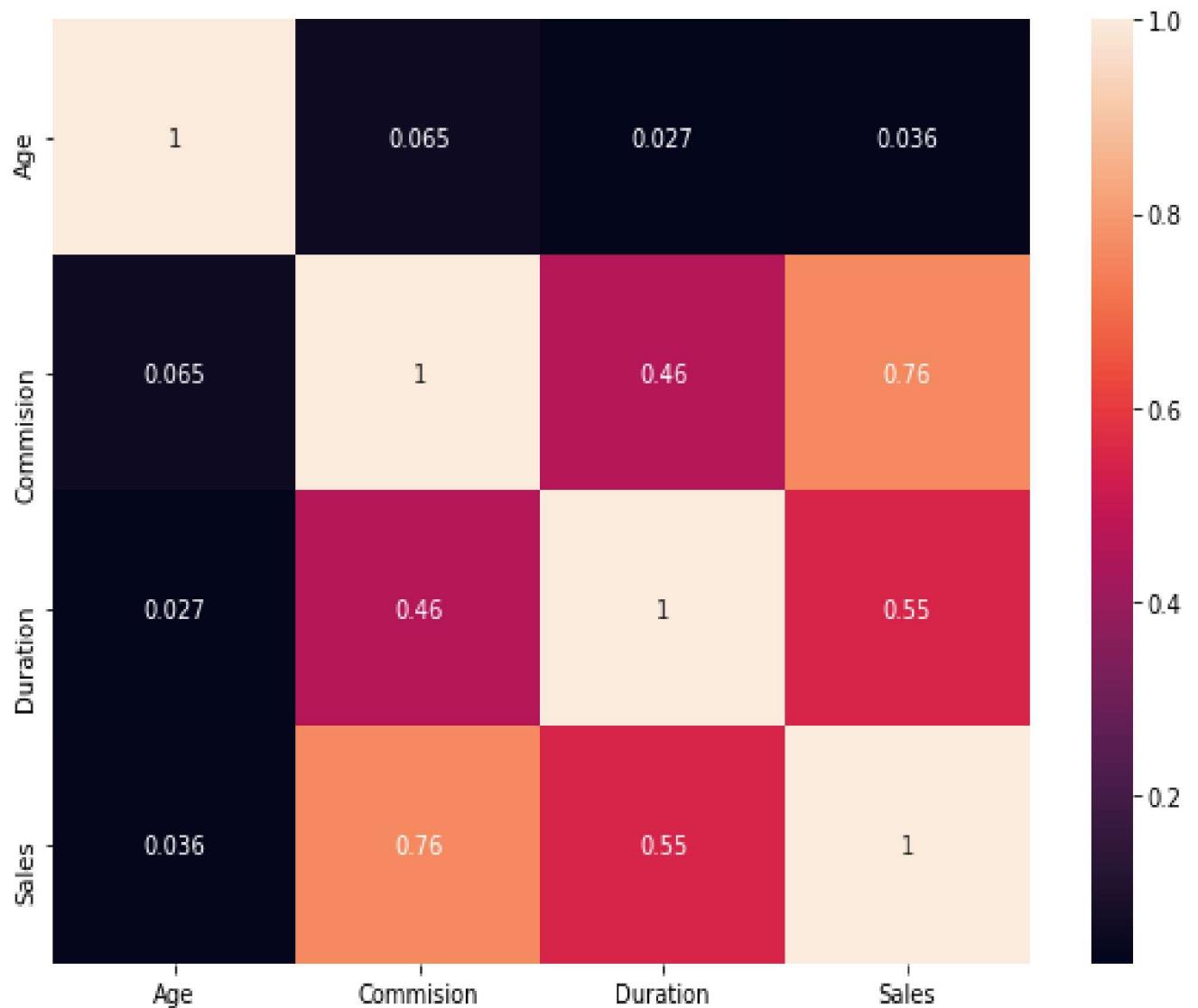
**Fig\_1.6**

**Conclusion | insight:** All columns data single behaviour show in fig\_1.5& fig\_1.6

## Bivariate Analysis

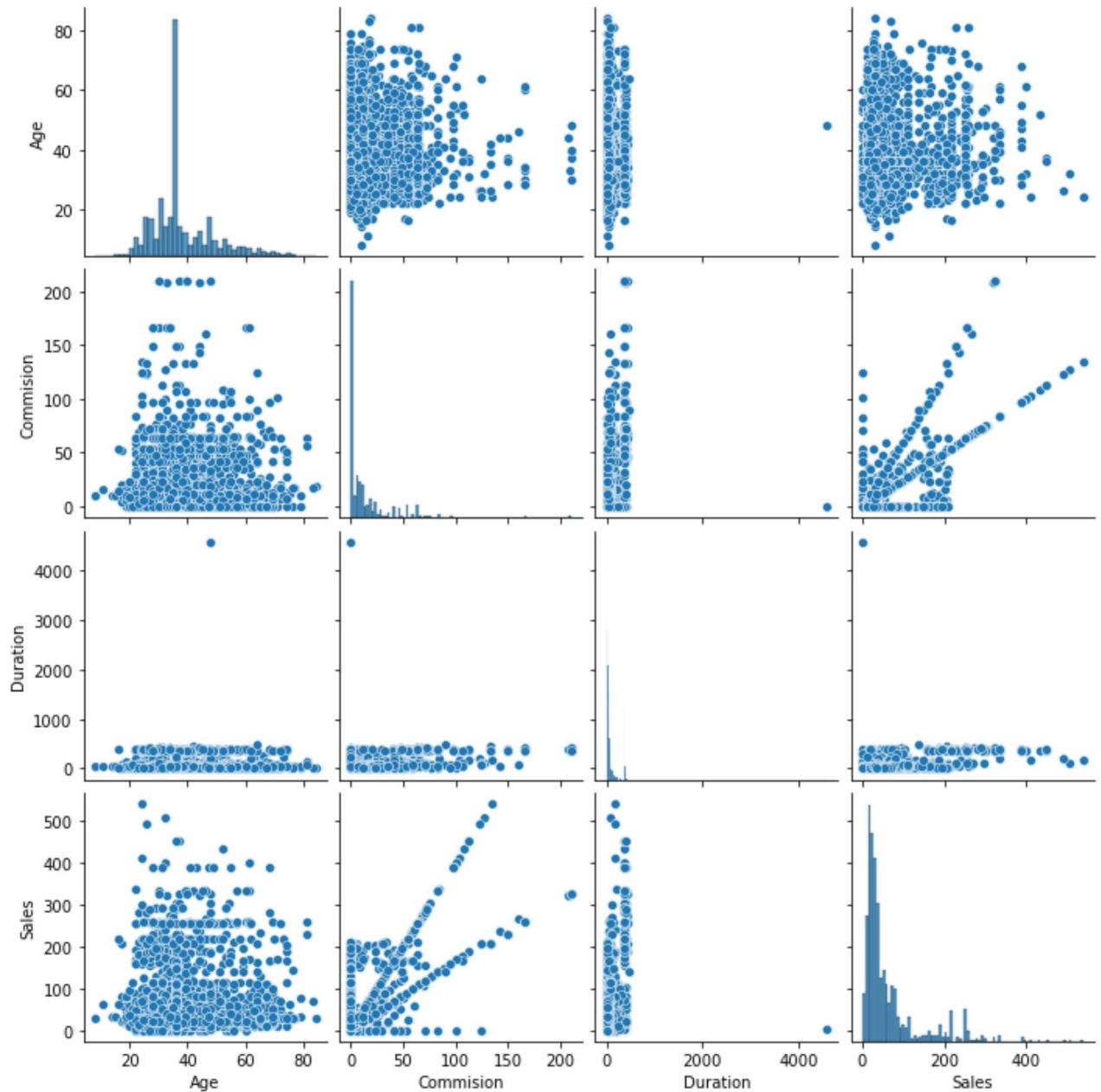
### Column- Claimed graph

**Fig\_1.7**

**Fig\_1.8**

**Conclusion | insight:** Correlation of all two data behaviour in show in fig\_1.8

## Multivariate Analysis



Fig\_1.9

**Conclusion | insight:** All data pair correction and behaviour show in fig\_1.9.

**Q 2.2 Data Split:** Split the data into test and train, build classification model CART, Random Forest, Artificial Neural Network.

**Solution –**

**Converted into categorical**

Categorised	Claimed	Type	Agency Code	Channel	Product	Destination
0	No'	Airlines	C2B	Online	Customised Plan	'ASIA'
1	Yes	Travel Agency	EPX	Offline	Cancellation Plan	'Americas'
2			CWT		Bronze Plan	EUROPE
3			JZI		Silver Plan'	
4					Gold Plan	

Table\_1.0

**Converted into categorical numerical data**

```
Int64Index: 2861 entries, 0 to 2999
Data columns (total 10 columns):
 #   Column      Non-Null Count  Dtype  
 --- 
 0   Age          2861 non-null    int64  
 1   Agency_Code  2861 non-null    int8   
 2   Type         2861 non-null    int8   
 3   Claimed      2861 non-null    int8  
 4   Commision    2861 non-null    float64 
 5   Channel      2861 non-null    int8  
 6   Duration     2861 non-null    int64  
 7   Sales         2861 non-null    float64 
 8   Product Name 2861 non-null    int8  
 9   Destination   2861 non-null    int8  
 dtypes: float64(2), int64(2), int8(6)

```

Fig\_1.10

## Categorical sample data

	Age	Agency_Code	Type	Claimed	Commision	Channel	Duration	Sales	Product Name	Destination
0	48	0	0	0	0.70	1	7	2.51	2	0
1	36	2	1	0	0.00	1	34	20.00	2	0
2	39	1	1	0	5.94	1	3	9.90	2	1
3	36	2	1	0	0.00	1	4	26.00	1	0
4	33	3	0	0	6.30	1	53	18.00	0	0

Fig\_1.11

## Classification model- CART train and test

### Grid search for CART

```
GridSearchCV(cv=3, estimator=DecisionTreeClassifier(),
            param_grid={'max_depth': [7, 8, 9, 10],
                        'min_samples_leaf': [15, 20, 25],
                        'min_samples_split': [45, 60, 75]})
```

Fig\_1.11b

Performing Grid Search and Cross validation so that best model performance can be obtained.

## Best Parameters

```
{'max_depth': 8, 'min_samples_leaf': 20, 'min_samples_split': 45}
```

Fig\_1.11c

Using these best parameters, we generate a Decision Tree and build a CART Model and find out the most important variable who will be most decisive for understanding the claim status.

## Variable Importance

	Imp
Agency_Code	0.473051
Sales	0.256875
Duration	0.089125
Commision	0.067888
Age	0.067663
Product Name	0.037309
Destination	0.008089
Type	0.000000
Channel	0.000000

Conclusion | insight: After grid search CART Modal data use all data is ready for ROC and AUC, matrix preformed.

## Random Forest- train and test

### Grid search for RF

```
GridSearchCV(cv=3, estimator=RandomForestClassifier(),
            param_grid={'max_depth': [7, 10], 'max_features': [4, 6],
                        'min_samples_leaf': [50, 100],
                        'min_samples_split': [150, 300],
                        'n_estimators': [301, 501]})
```

Performing Grid Search and Cross validation so that best model performance can be obtained.

## Best Parameters

```
{'max_depth': 7,
'max_features': 4,
'min_samples_leaf': 50,
'min_samples_split': 150,
'n_estimators': 301}
```

Min max estimators using in RF Model



## Best Parameters

```
{'activation': 'relu',
 'hidden_layer_sizes': (100, 100, 100),
 'max_iter': 10000,
 'solver': 'adam',
 'tol': 0.01}
```

**Conclusion | insight:** After grid search ANN Modal data use all data is ready for ROC and AUC, matrix preformed.

**Q 2.3 Performance Metrics:** Comment and Check the performance of Predictions on Train and Test sets using Accuracy, Confusion Matrix, and Plot ROC curve and get ROC\_AUC score, classification reports for each model.

### Solution –

#### 1. CART Model for Train -

**Confusion matrix –**

```
[1199, 160]
[ 262, 381]
```









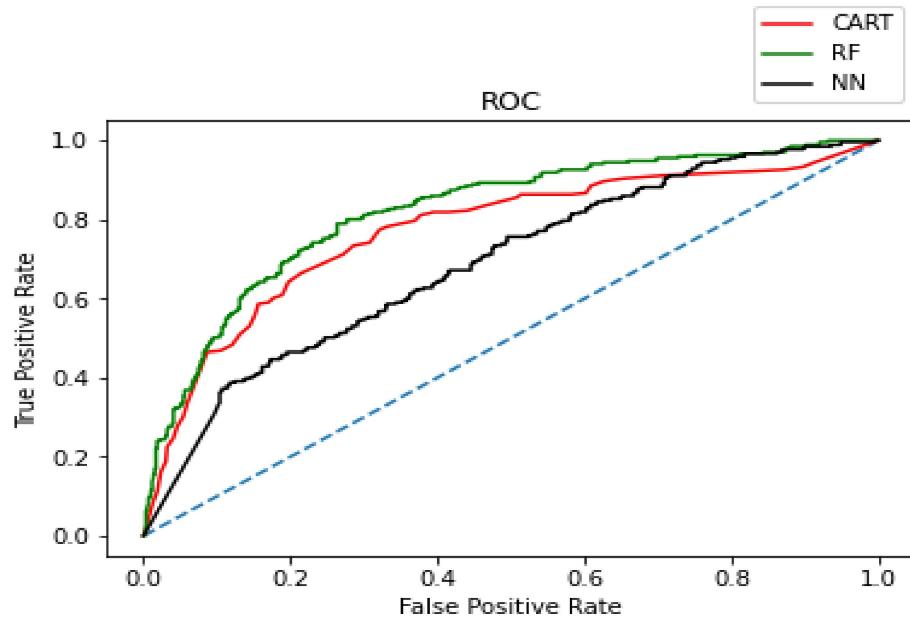






Fig\_1.18

## ROC Curve for the CART, RF and ANN Test



Fig\_1.19

**Conclusion | insight:** As per compare all model RF Model is the best and CART model performs are also good but ANN model performs is low but all model are best performs.

**Q 2.5 Inference:** Based on the whole Analysis, what are the business insights and recommendations.

## **Solution –**

**Conclusion | insight:**

1. This is understood by looking at the insurance data by drawing relations between different variables such as day of the incident, time, age group, and associating it with other external information such as location, behaviour patterns, weather information, airline/vehicle etc.
2. Almost need to train the JZI agency resources to pick up sales as they are in bottom, need to run promotional marketing campaign or evaluate if we need to tie up with alternate agency.
3. As per the data 90% of insurance is done by online channel. Other interesting fact is more sales happen via Agency than Airlines and the trend shows the claim are processed more at Airline.
4. Reduce claim handling costs Insights gained from data and AI-powered analytics could expand the boundaries of insurability, extend existing products, and give rise to new risk transfer solutions in areas like a non-damage business interruption and reputational damage.