

Historical Data Losses

Full Project Report

- DIPANSHU RANGA
-dipanshuranga@gmail.com

Historical Losses Dataset

Analysis Report

ABOUT:

Analysis Report on the past data of **losses** from an **insurance company** on **automobile** insurance policy, to **predict future losses** for same featured data as this dataset to decide **amount of premium of new customer's**.

1. DATA CONSOLIDATION:

Original dataset named with **“Historical Losses Data.csv”**.

Dataset comprises of 15,290 observations and 9 variables

1st variable is **“Policy Number”** i.e., Not-**Significant** for making predictions so it’s **dropped**.

Now, Dataset has **15,290 observations** and **8 variables**.

2. EXPLORATORY DATA ANALYSIS:

2.1. VARIABLE - IDENTIFICATION:

Out of 8 variables there is **one dependent** variable and other **seven** are **independent** variables.

Dependent Variable: --

- 1. Losses

Independent Variables: --

- 1. Age
- 2. Years_of_Driving_Experience
- 3. Number_of_Vehicles
- 4. Gender
- 5. Married
- 6. Vehicle_Age
- 7. Fuel_Type

*There are no NULL Values in dataset.

Data Type:	
Numeric:	Character:
1. Age	1. Gender
2. Years_of_Driving_Experience	2. Married
3. Number_of_Vehicles	3. Fuel_Type
4. Vehicle_Age	
5. Losses	

*All Numeric variables are integers.

Variable Category:	
Continuous:	Categorical:
1. Age	1. Gender
2. Years_of_Driving_Experience	2. Married
3. Number_of_Vehicles	3. Fuel_Type
4. Vehicle_Age	
5. Losses	

2.2. UNIVARIATE ANALYSIS:

Range Information:

Variable	Range/Category:
Age	16 – 70
Years_of_Driving_Experience	0 - 53
Number_of_Vehicles	1 - 4
Gender	‘M’ and ‘F’
Married	‘Married’ and ‘Single’
Vehicle_Age	0 – 15
Fuel_Type	‘P’ and ‘D’
Losses	13 - 3500

*‘M’ and ‘F’ stands for ‘Male’ and ‘Female’ respectively. **‘P’ and ‘D’ stands for ‘Petrol’ and ‘Diesel’ respectively.

Numerical-Data Variable Description:

	Age	Years_of_Driving_Experience	Number_of_Vehicles	Vehicle_Age	Losses
Count	15290	15290	15290	15290	15290
Mean	42.32	23.73	2.49	8.65	389.85
Std	18.28	17.85	0.95	4.34	253.72
Min	16.0	0.0	1.0	0.0	13.0
25%	24.0	6.0	2.0	6.0	226.0
50%	42.0	23.0	2.0	9.0	355.0
75%	61.0	42.0	3.0	12.0	489.0
Max	70.0	53.0	4.0	15.0	3500.0

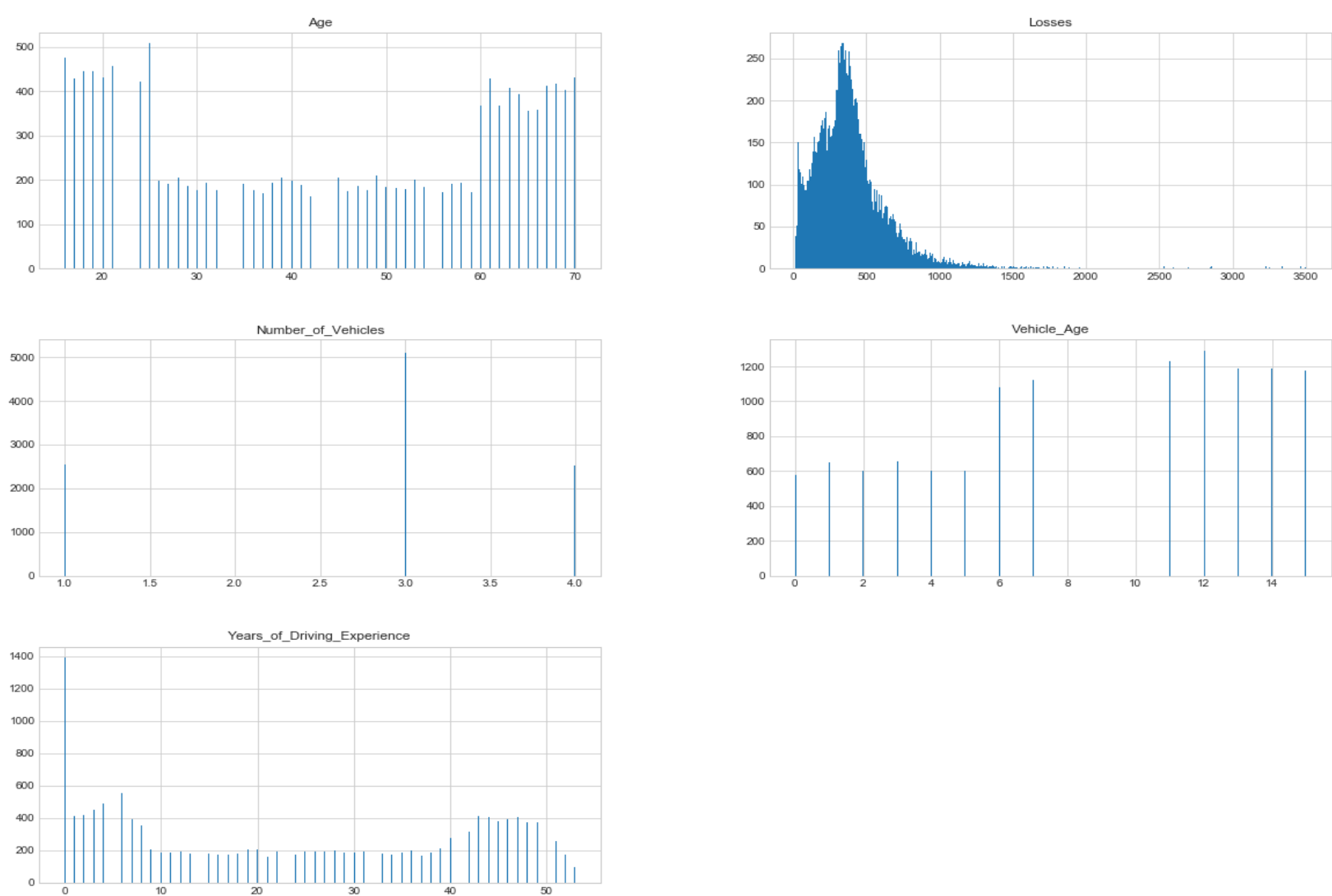
- There are no missing values in numerical columns.
- Maximum observation of “Losses” is 3500 and 75 percentile is 489 which tell us that “**Losses**” contain **outlier** values.

Skewness:

Variable	Skewness
Age	+0.05
Years_of_Driving_Experience	+0.097
Number_of_Vehicles	+0.0065
Vehicle_Age	-0.344
Losses	+2.55

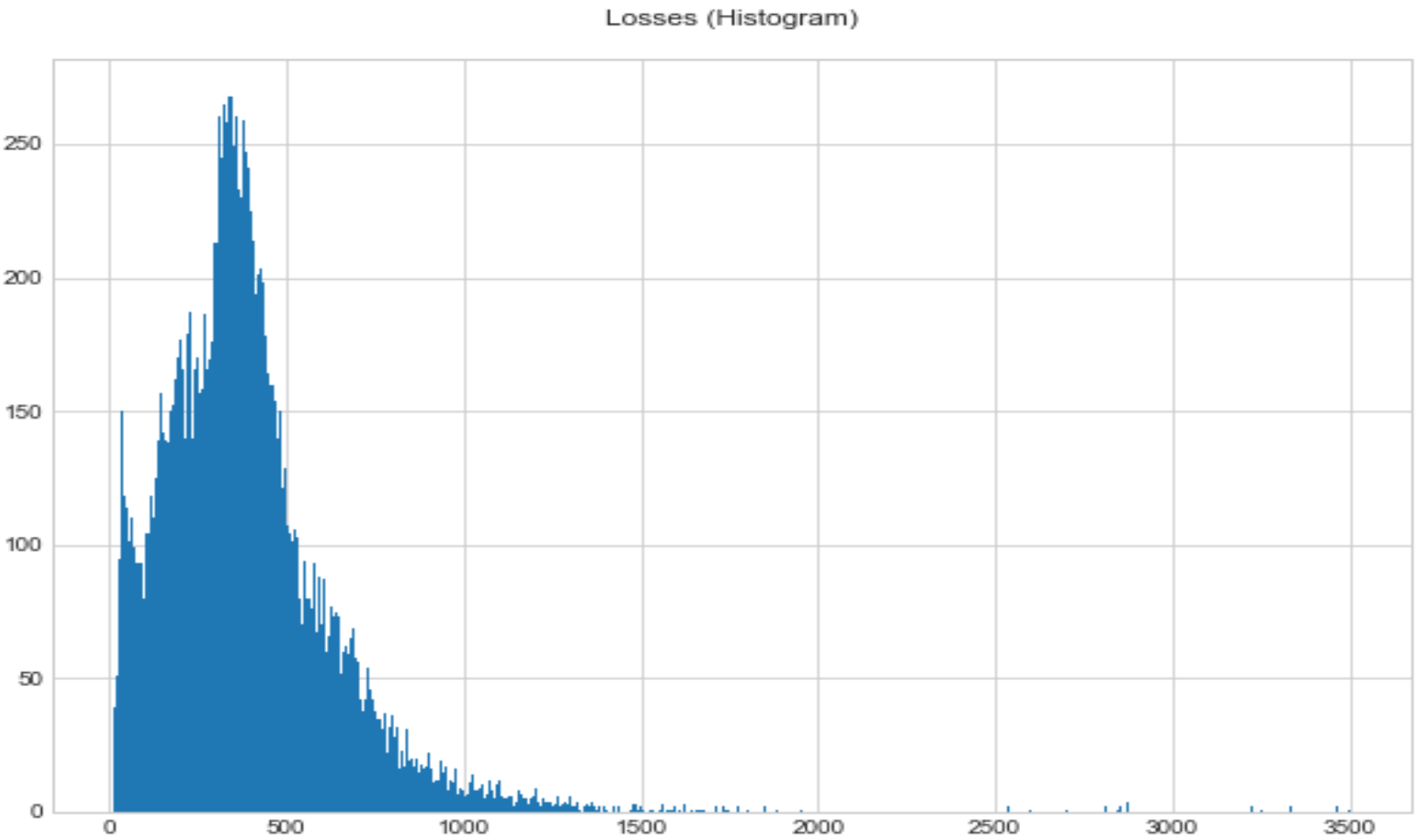
- “Losses” is **highly positive skewed** with a skewness of over +2.0
- “Vehicle_Age” is **negatively skewed** with skewness of over -0.3.

Histograms:



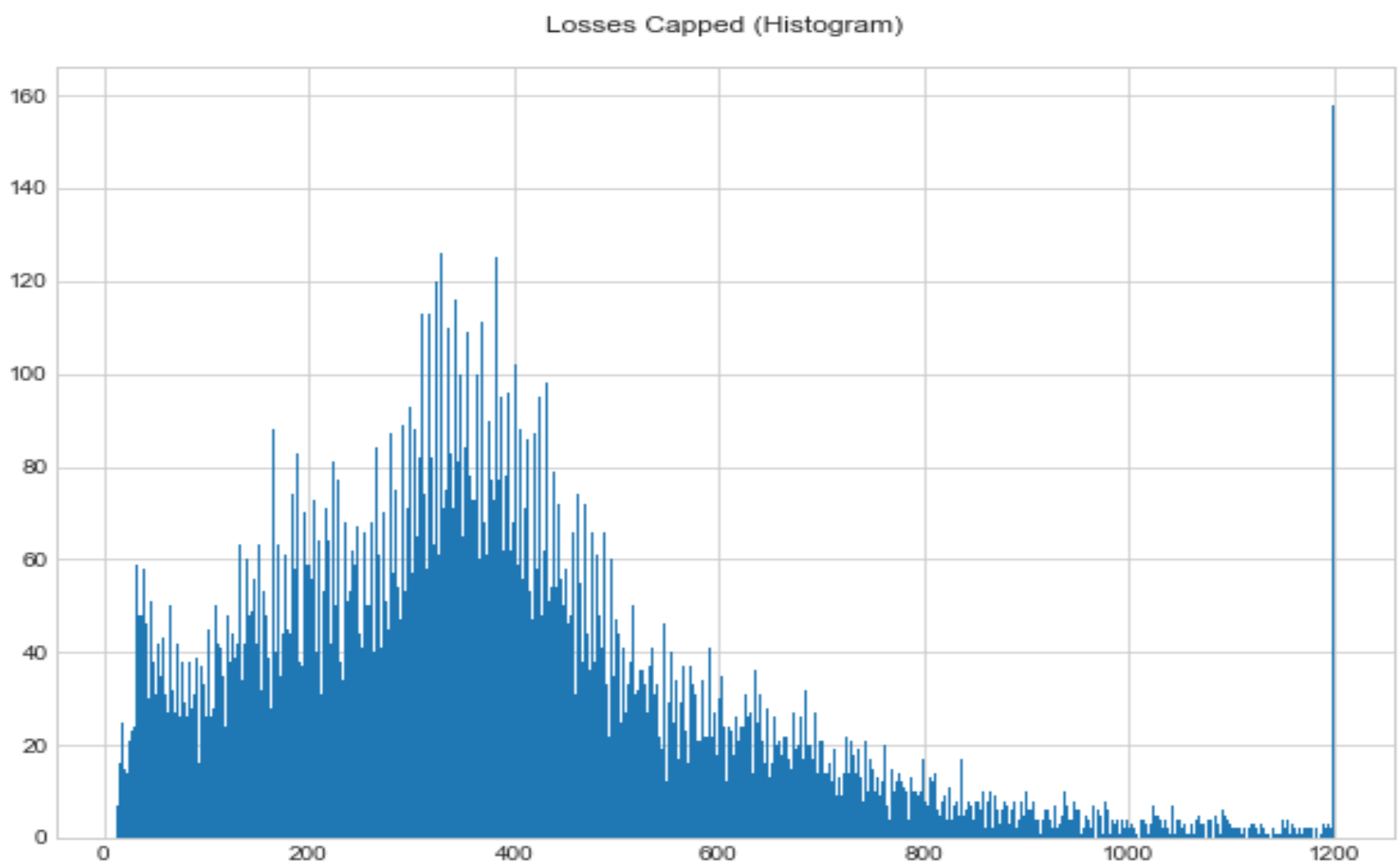
- Numbers of population in age group of ‘30 to 60’ are **lower** in **number** in comparison to other age group.
- People with **zero** driving **experience** are **high** in **number**.

- Taking a separate look at Losses Histogram:
-



- Already observed high skewness of 2+ in Skewness-Table.
- There are people with **high losses**, but they are **low** in **number**.
- We can put all the **high losses** observations in **one category**.

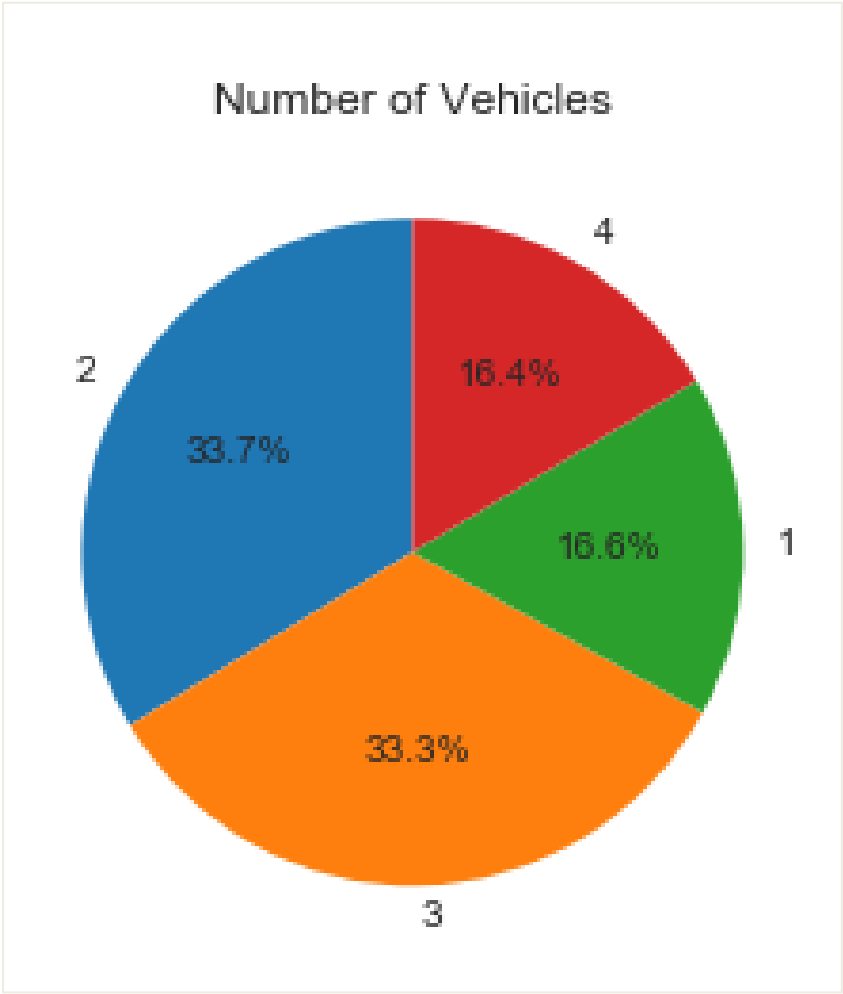
- After upper capping of “Losses” at 1200:



- The skewness is gradually decreased to 1.05 (before capping it was 2+).

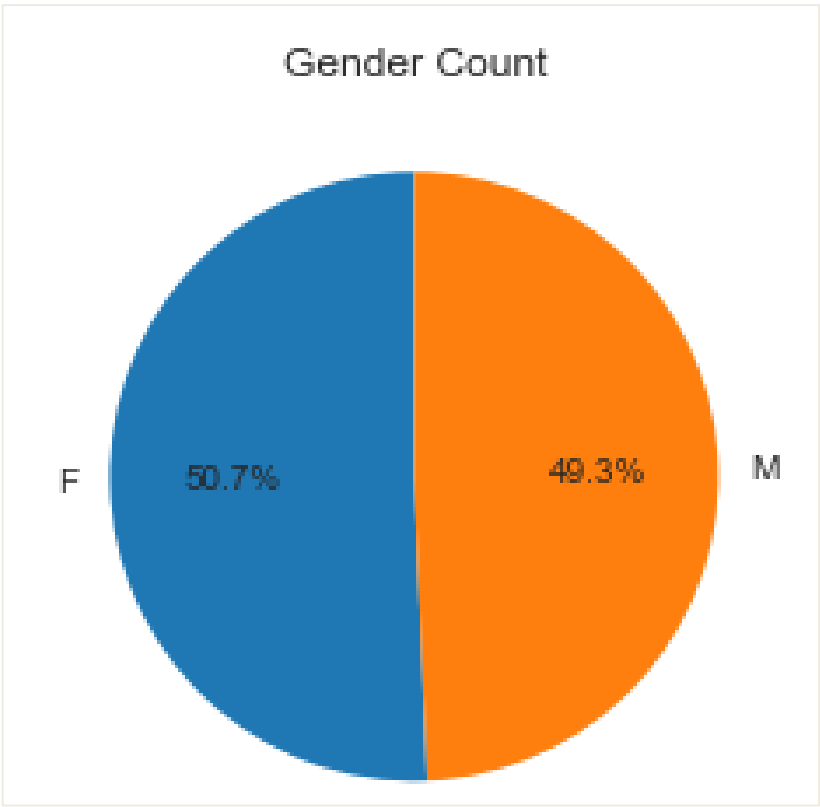
Pie-Plots:

- “Number of Vehicles”:

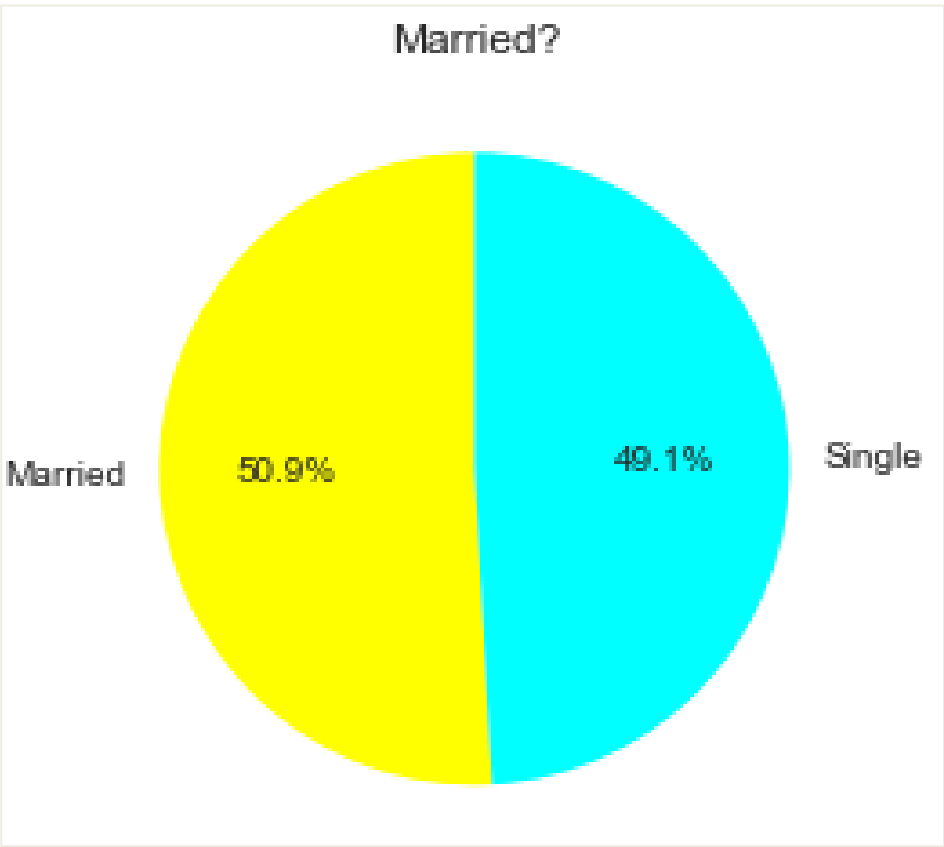


- People owning ‘two’ and ‘three’ vehicles are nearly **double** in number who owns ‘one’ or ‘four’ vehicles.

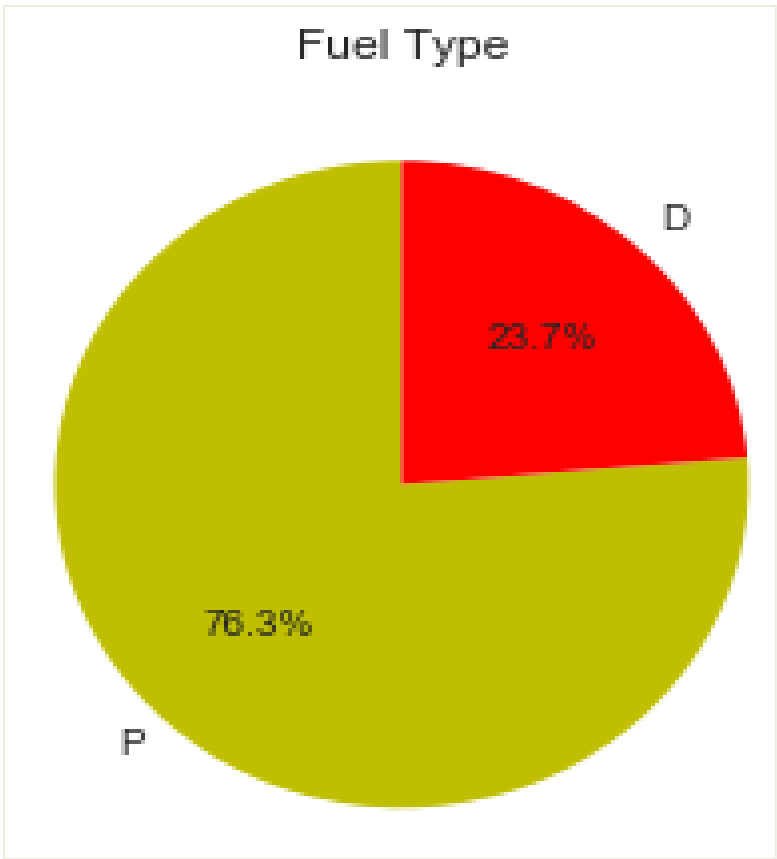
- Categorical Data (Gender, Married & Fuel_Type):



■ **Male & Female** have nearly equal share in observations.



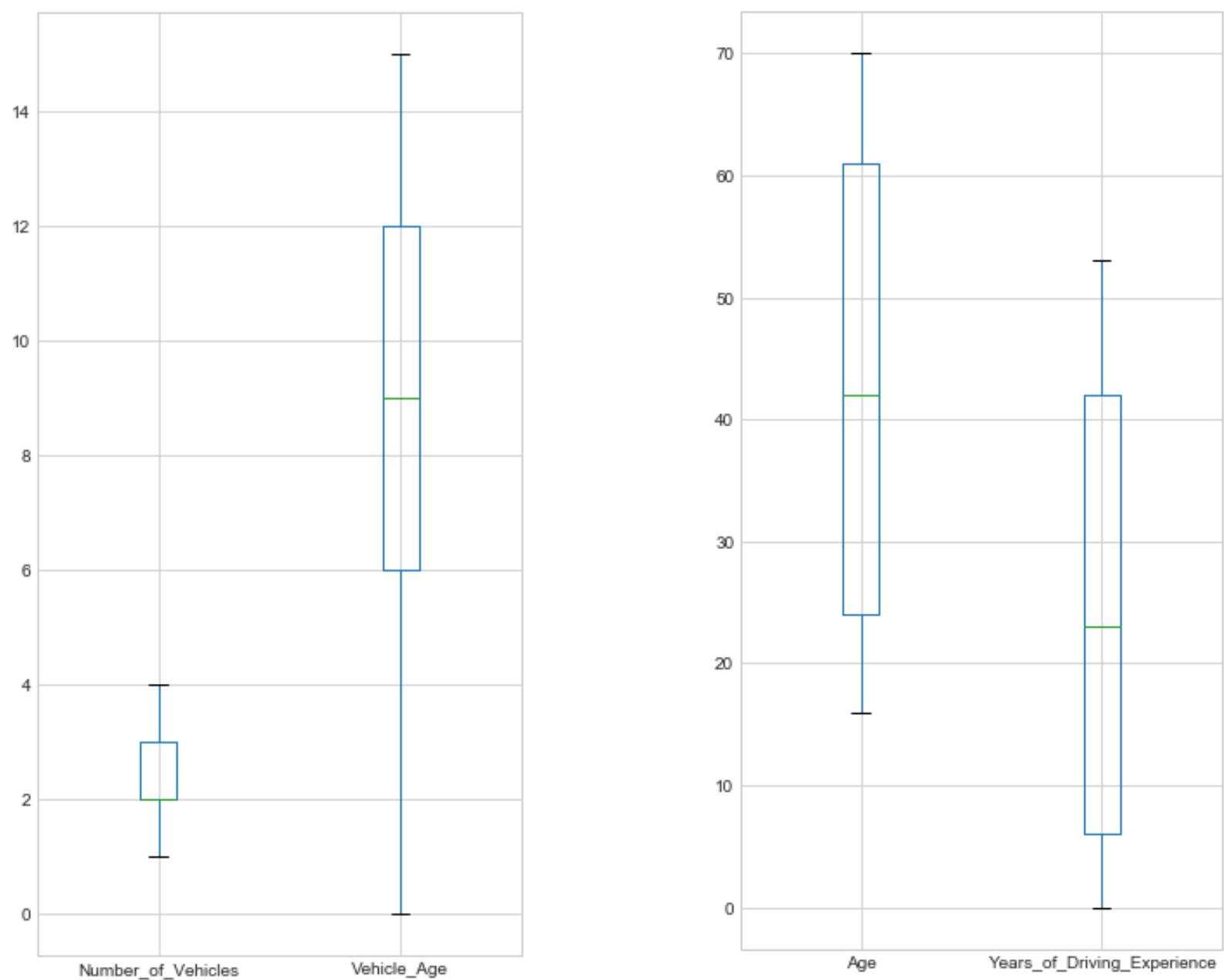
■ **Married & Single** population both have nearly equal number of observation.



■ **Petrol** type vehicles are **more** in number.

■ **Petrol** Type vehicles are nearly **3x** of **Diesel** Type Vehicles.

Box-Plots:



■ There are No Outliers in Age, Years_of_Driving_Experience, Number_of_Vehicles and Vehicle_Age.

■ “Age” & “Years_of_Driving_Experience” has nearly same scale.

- **Univariate view for Bivariate Analysis:**

As seen in **Univariate Analysis** only “Gender” & “Married” have nearly **same number** of **observation** for every **unique value**. But other variables have different number of observations for different unique values.

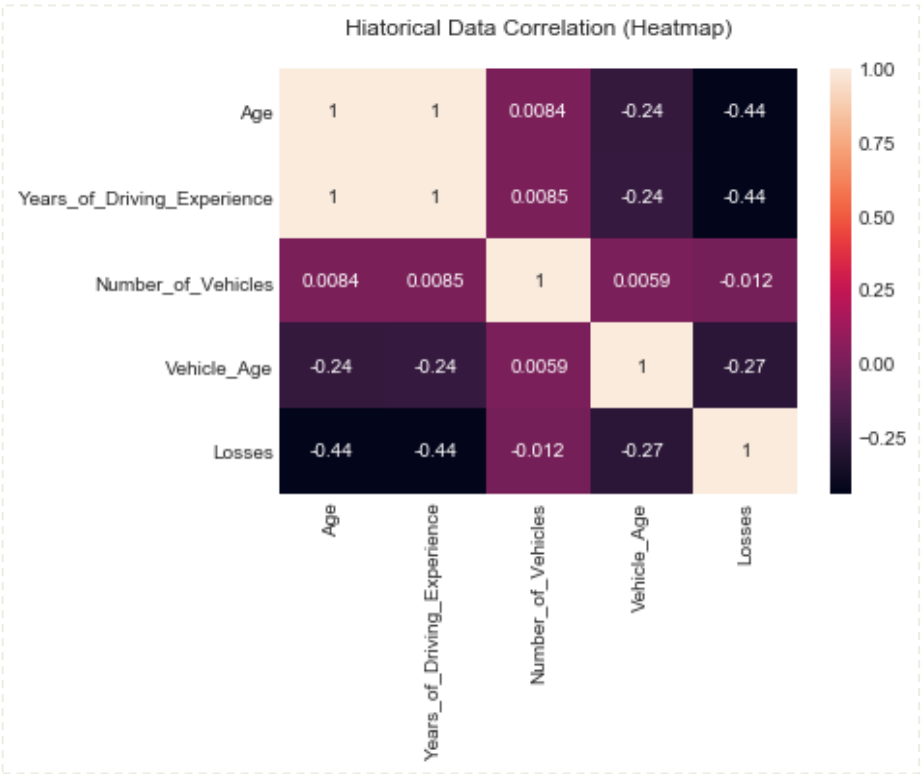
For Example:

1. There are more number of people who own two & three vehicles then who own one & four
2. People who own Petrol Vehicle are three times than the people who own Diesel Vehicle.

So, for **Bivariate Analysis** we are **taking plot of Average of Losses** (dependent variable) with every other variable (**instead of sum of losses**).

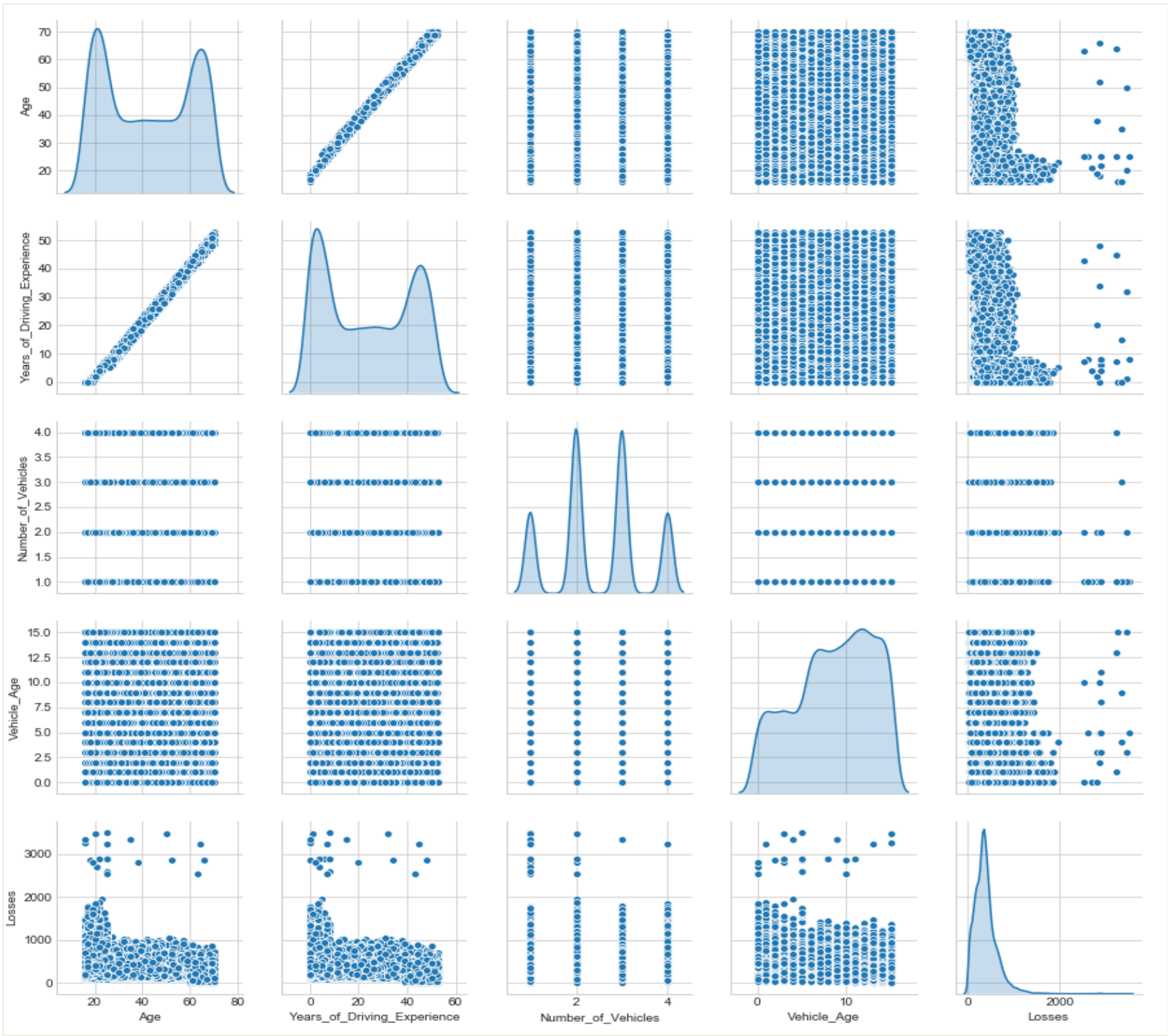
2.3. BIVARIATE ANALYSIS:

Correlation Heat-Map:



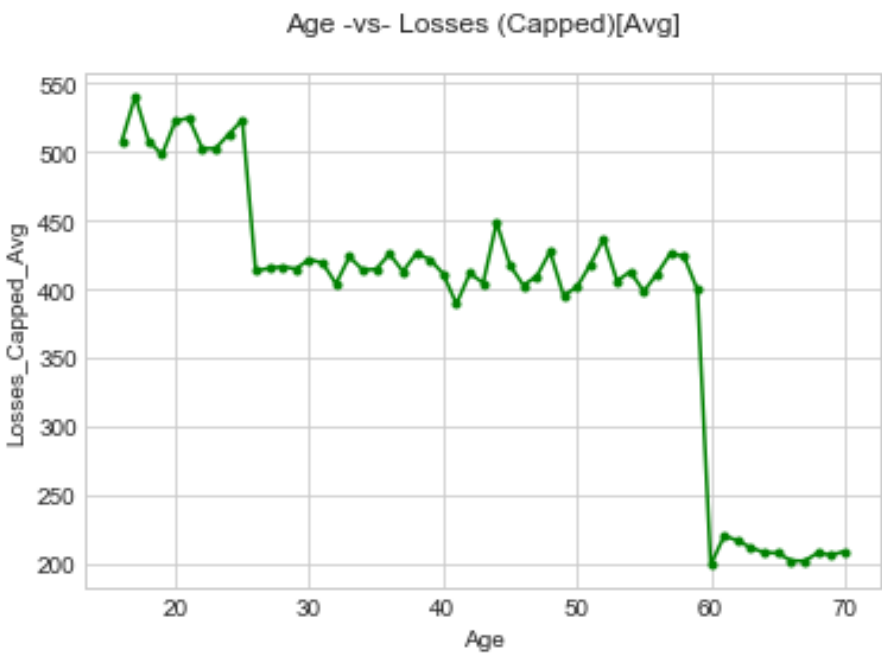
■ There is **high positive correlation** between “Age” and “Years_of_Driving_Experience”

Pair-Plot:



- There are observations in population with **high losses**, but they are very **low in number**.
- It also **confirms** the **high correlation** between “Age” and “Years_of_Driving_Experience”
- The skewness of “**Vehicle_Age**” is natural because **most** number of **populations** owns **old vehicles**.

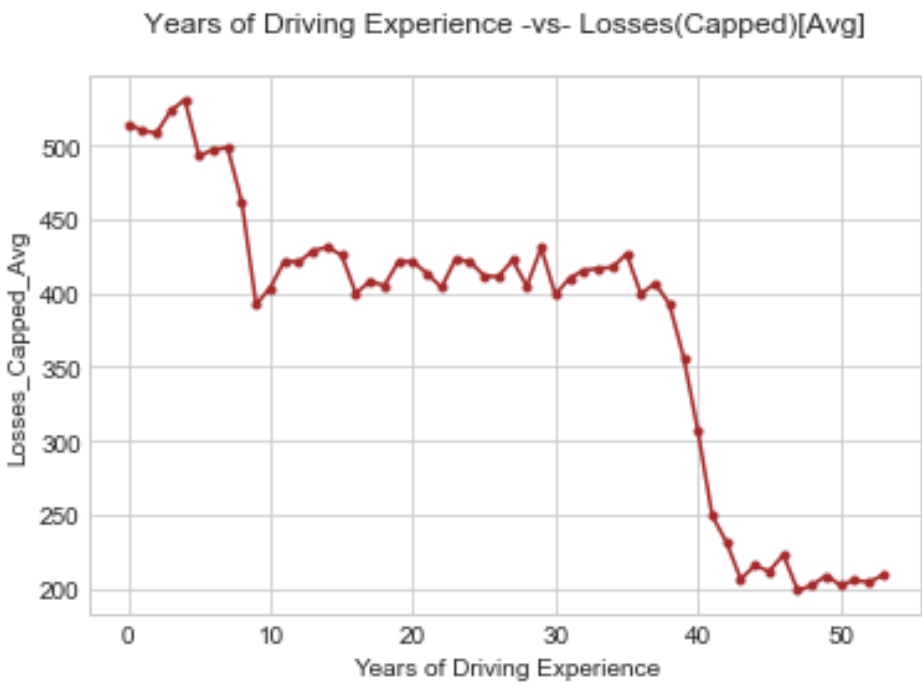
2.3.1. Age vs. Losses:



- People with young **Age group** like '16_to_25' have **very high losses**.
- People with **Age group** like '26_to_59' have **moderate high losses**.
- People with young **Age group** like '60_to_70' have **low losses**.

*Bucketing Needed.

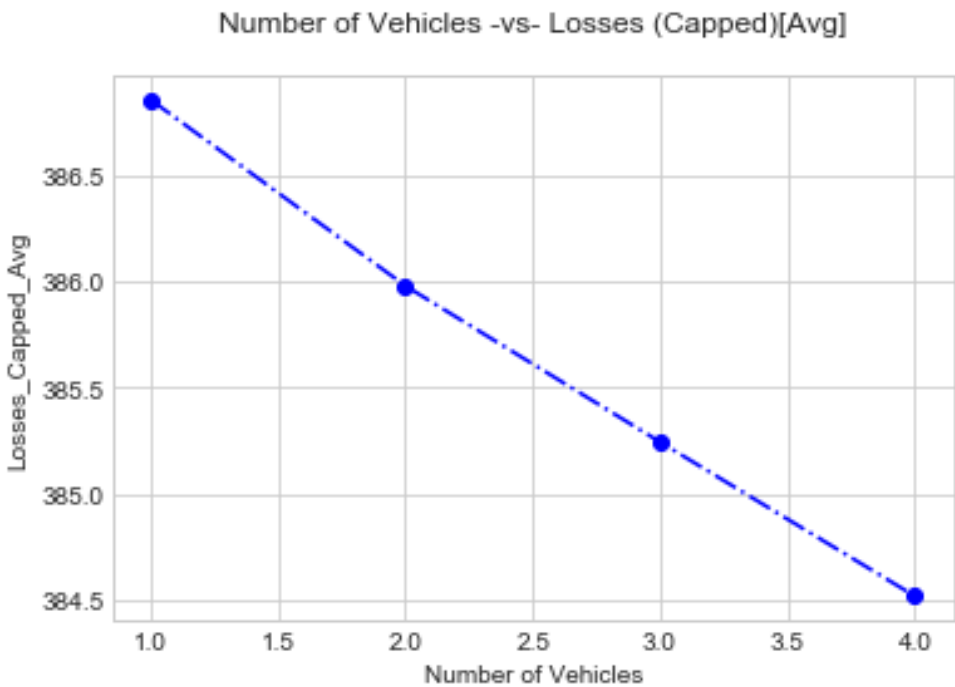
2.3.2. Years_of_Driving_Experience vs. Losses:



- People with low Driving **Experience** like **0 to 8** have **very high losses**.
- People with Driving **Experience** of **9 to 40** have **moderate losses**.
- People with Driving **Experience** of **41 to 53** have **low losses**.

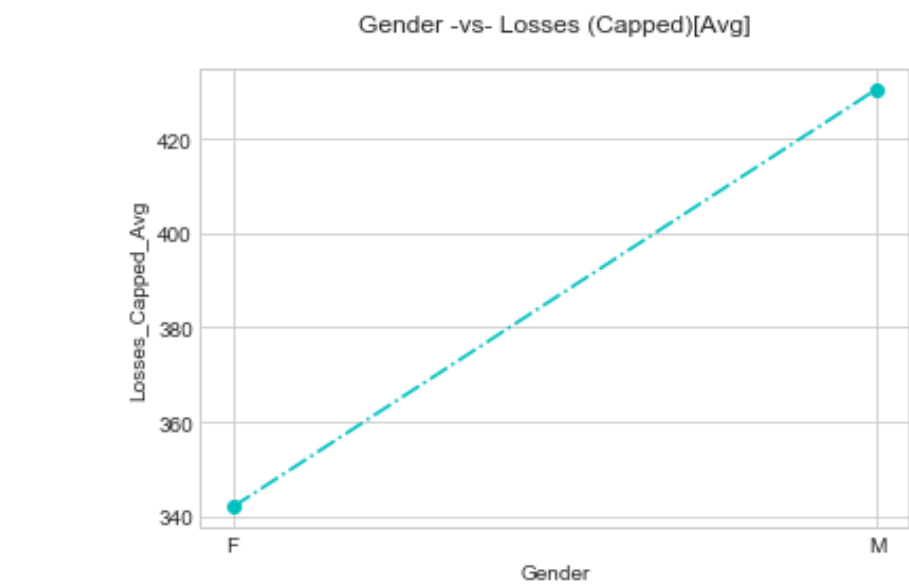
*Bucketing Needed.

2.3.3. Number_of_Vehicles vs. Losses:



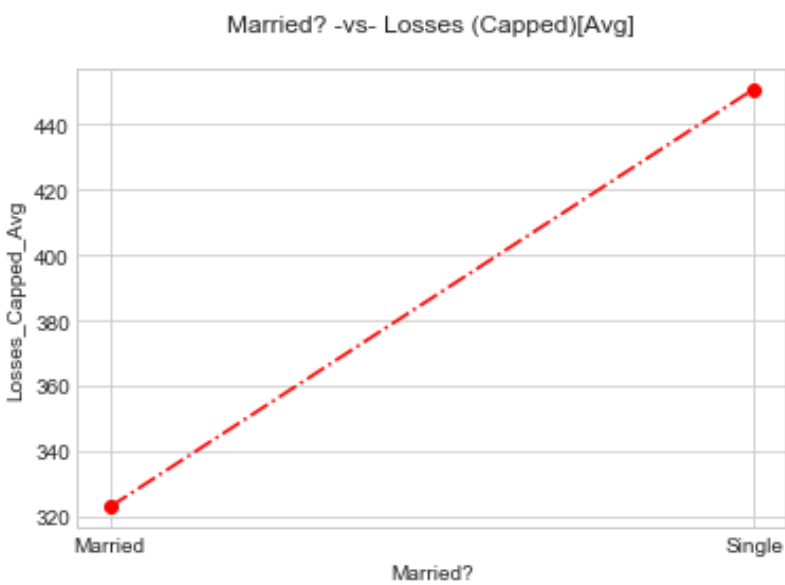
- People with **one vehicle** have more **losses** than others.

2.3.4. Gender vs. Losses:

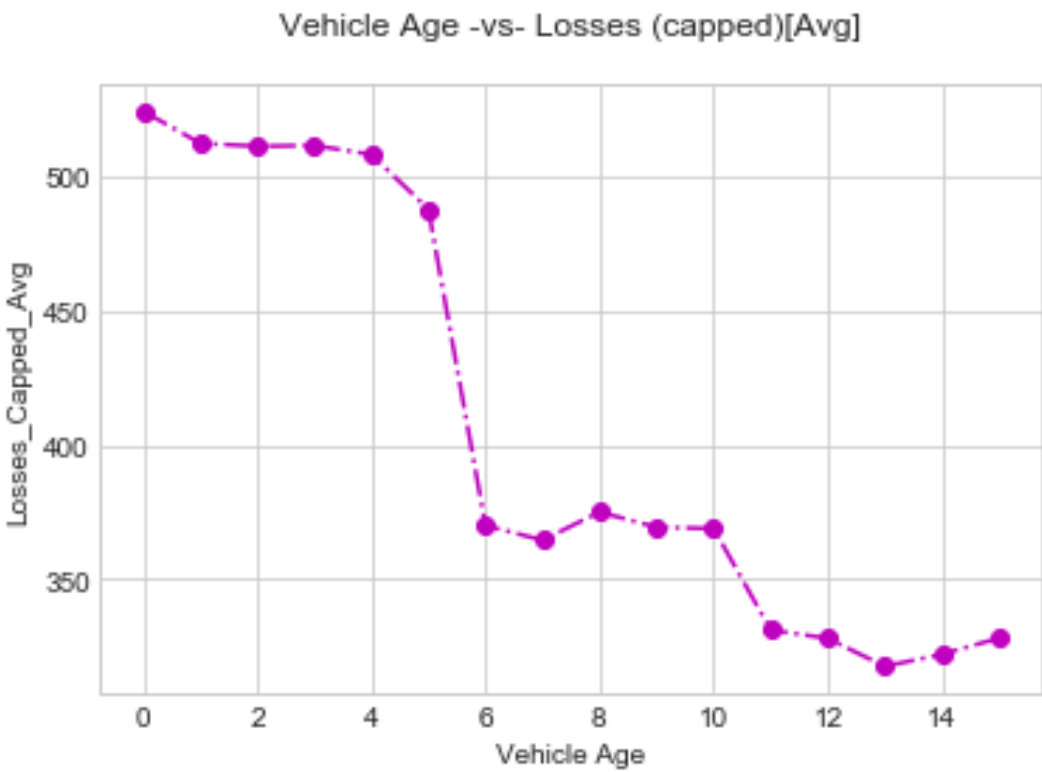


- Male population has **higher** number of **losses**.
- Single people have **higher** number of **losses**.

2.2.5. Married vs. Losses:



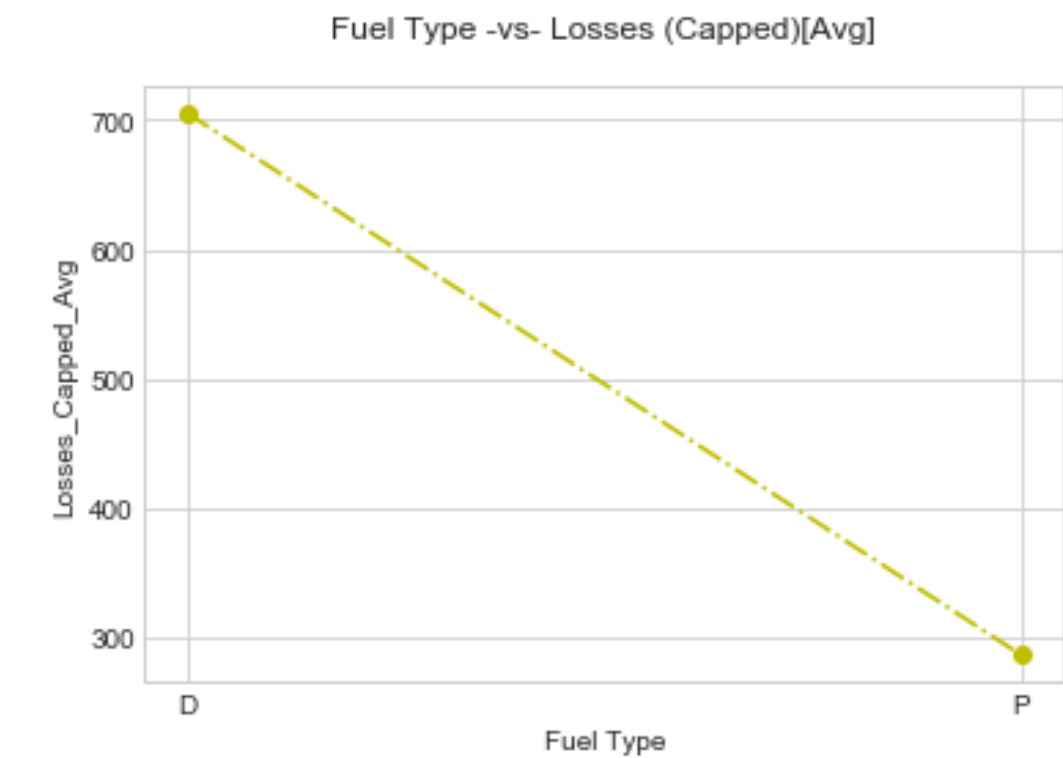
2.2.6. Vehicle_Age vs. Losses:



- People who own **new vehicles (0 to 5 years old)** have **high losses**.
- People who own **vehicles** aging from ‘**6 to 10**’ has comparably **low losses**. **But**, people with vehicle age group of ‘**11 to 15**’ have **even more lower losses**.

*Bucketing Needed.

2.2.7. Fuel_Type vs. Losses:



- People with **Diesel** Type vehicles have **high** amount of **losses**.

- **EDA Points Summary:**
 - There are no missing values in numerical columns.
 - Maximum observation of “Losses” is 3500 and 75 percentile is 489 which tell us that “**Losses**” contain **outlier** values.
 - “Losses” is **highly positive skewed** with a skewness of over +2.0
 - “Vehicle_Age” is **negatively skewed** with skewness of over -0.3.
 - Numbers of population in age group of ‘**30 to 60**’ are **lower** in **number** in comparison to other age group.
 - People with **zero** driving **experience** are **high** in **number**.
 - Already observed high skewness of 2+ in Skewness-Table.
 - There are people with **high losses**, but they are **low** in **number**.
 - We can put all the **high losses** observations in **one category**.
 - The skewness is gradually decreased to 1.05 (before capping it was 2+).
 - People owning ‘two’ and ‘three’ vehicles are nearly **double** in number who owns ‘one’ or ‘four’ vehicles.
 - **Male & Female** have nearly equal share in observations.
 - **Married & Single** population both have nearly equal number of observation.
 - **Petrol** type vehicles are **more** in number.
 - **Petrol** Type vehicles are nearly **3x** of **Diesel** Type Vehicles.
 - There are **no outliers** in **Age**, **Years_of_Driving_Experience**, **Number_of_Vehicles** and **Vehicle_Age**.
 - **Age & Years_of_Driving_Experience** has nearly **same scale**.
 - There is **high positive correlation** between “Age” and “Years_of_Driving_Experience”
 - There are observations in population with **higher losses**, but they are very **low in number**.
 - It(pair-plot) also **confirms** the **high correlation** between “Age” and “Years_of_Driving_Experience”
 - The skewness of “**Vehicle_Age**” is natural because **most** number of **populations** owns **old vehicles**.
 - People with young **Age group** like ‘**16 to 25**’ have **very high losses**.
 - People with **Age group** like ‘**26 to 59**’ have **moderate high losses**.
 - People with young **Age group** like ‘**60 to 70**’ have **low losses**.
 - People with low Driving **Experience** like **0 to 8** have **very high losses**.
 - People with Driving **Experience** of **9 to 40** have **moderate losses**.
 - People with Driving **Experience** of **41 to 53** have **low losses**.
 - People with **one vehicle** have more **losses** than others.
 - **Male** population has **higher** number of **losses**.
 - **Single** people have **higher** number of **losses**.
 - People who own **new vehicles (0 to 5 years old)** have **high losses**.
 - People who own **vehicles** aging from ‘**6 to 10**’ has comparably **low losses**. **But**, people with vehicle age group of ‘**11 to 15**’ have **even more lower losses**.
 - People with **Diesel** Type vehicles have **high** amount of **losses**.

2.3. Conclusion:

- From the above **EDA Report** we get to know that “**Age**”, “**Years_of_Driving_Experience**” and “**Vehicle_Age**” needs **bucketing**.

2.4. Bucketing:

2.4.1. Age Bucketing:

Age Group **16 to 25** as → **21**
Age Group **26 to 59** as → **43**
Age Group **60 to 70** as → **65**

2.4.2. Years_of_Driving_Experience Bucketing:

YODE Group **0 to 8** as → **4**
YODE Group **9 to 40** as → **25**
YODE Group **41 to 53** as → **47**

2.4.3. Vehicle_Age Bucketing:

Vehicle_Age **0 to 5** as → **3**
Vehicle_Age **6 to 10** as → **8**
Vehicle_Age **11 to 15** as → **13**

3. Confirmatory Data Analysis:

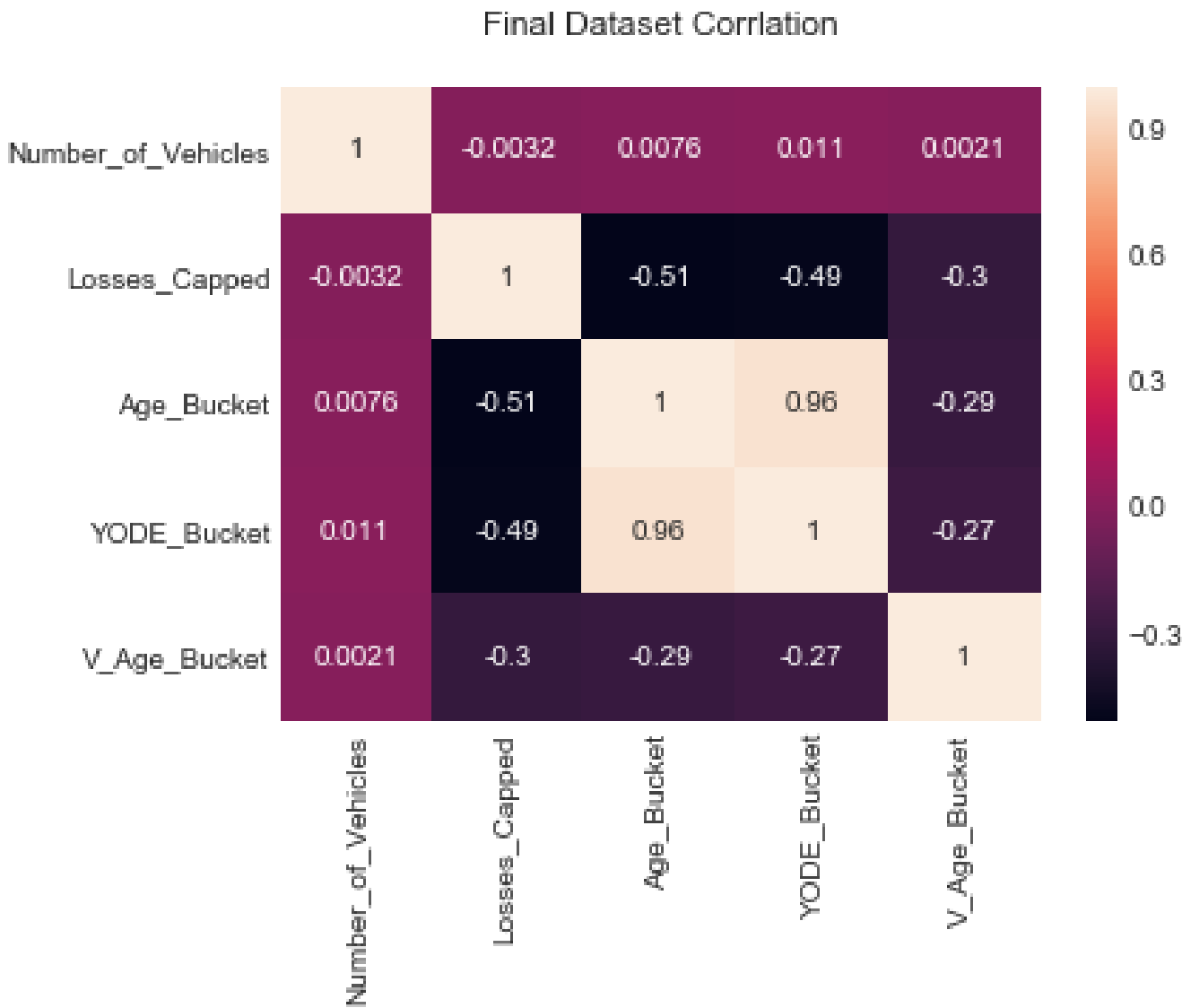
3.1. Data Information:

	Number_of_Vehicles	Losses_Capped	Age_Bucket	YODE_Bucket	V_Age_Bucket
count	15290.0	15290.0	15290.0	15290.0	15290.0
mean	2.49	385.63	42.68	24.266	8.781
std	0.953	228.77	16.78	16.61	3.9
min	1.0	13.0	21.0	4.0	3.0
25%	2.0	226.0	21.0	4.0	8.0
50%	2.0	355.0	43.0	25.0	8.0
75%	3.0	489.0	65.0	47.0	13.0
max	4.0	1200.0	65.0	47.0	13.0

Columns	Skewness
Number_of_Vehicles	0.0065
Losses_Capped	1.055
Age_Bucket	0.024
YODE_Bucket	0.118
V_Age_Bucket	-0.281

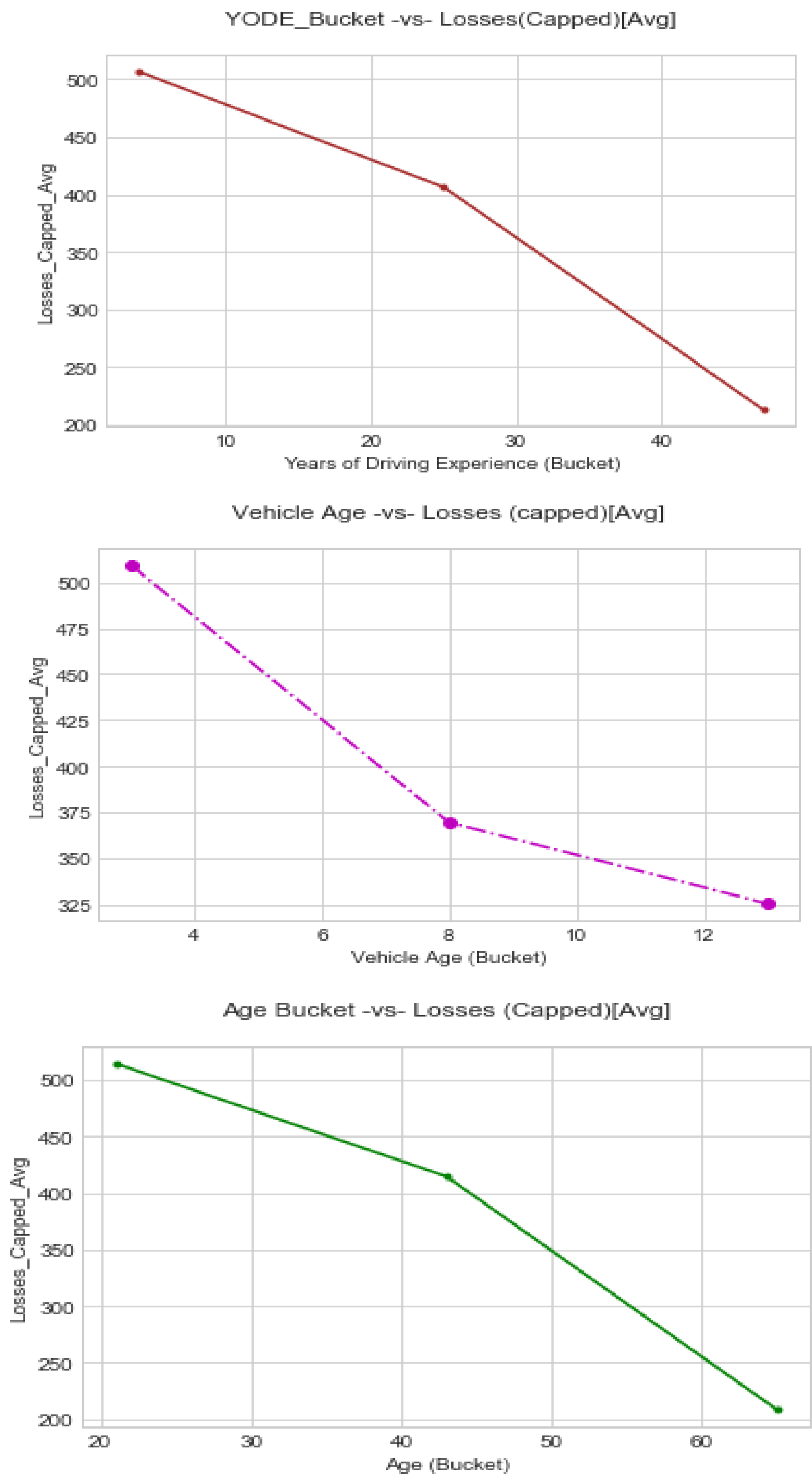
3.2. Bivariate(CDA):

.2.1 Correlation:



- As you can see “Age_Bucket” & “YODE_Bucket” still have high correlation. So, we will make two separate models in model preparation.
- YODE is used here and further in report as an abbreviation of Years_of_Driving_Experience.

3.2.2. Bucketed datasets Plots:



■ All three graphs showing **Linear Relationship** and a negative line.

4. Data Pre-Processing:

4.1. Missing Value Treatment:

- There are no missing values in this dataset.

4.2. Oulier Treatment:

- Outliers have been treated above with **capping** & **bucketing**.

4.3. Variable Transformation:

- Converting categorical data into numerical data.

4.4. Variable Creation:

- Creating two different data because Age and Years of Driving Experience still have correlation.

5. Model Development:

5.1. Model-1:

After the OLS through Backward Elimination on “**X_Age**” (the dataset contains ‘Age_Bucketed’ and others but not ‘YODE_Bucketed’ we got these results with a significance level of 5% (0.05)

OLS Regression Results						
Dep. Variable:	Losses_Capped	R-squared:	0.775			
Model:	OLS	Adj. R-squared:	0.775			
Method:	Least Squares	F-statistic:	8784.			
Date:	Fri, 15 Mar 2019	Prob (F-statistic):	0.00			
Time:	05:56:56	Log-Likelihood:	-1.0364e+05			
No. Observations:	15290	AIC:	2.073e+05			
Df Residuals:	15284	BIC:	2.073e+05			
Df Model:	6					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Number_of_Vehicles	95.8454	1.552	61.754	0.000	92.803	98.888
Gender	119.5762	3.450	34.655	0.000	112.813	126.340
Married	144.0847	3.513	41.010	0.000	137.198	150.971
Fuel_Type	-269.0490	5.092	-52.837	0.000	-279.030	-259.068
Age_Bucket	1.6484	0.102	16.235	0.000	1.449	1.847
V_Age_Bucket	13.2384	0.435	30.467	0.000	12.387	14.090
Omnibus:	185.030	Durbin-Watson:	1.180			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	265.037			
Skew:	0.148	Prob(JB):	2.81e-58			
Kurtosis:	3.573	Cond. No.	144.			

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

- $R^2 = 0.775$
- AIC = 2.073
- Durbin Watson = 1.80
- “Number_of_Vehicles” is eliminated as its p-values > our significance level

5.2. Model-2:

OLS Regression Results						
Dep. Variable:	Losses_Capped	R-squared:	0.772			
Model:	OLS	Adj. R-squared:	0.772			
Method:	Least Squares	F-statistic:	8644.			
Date:	Fri, 15 Mar 2019	Prob (F-statistic):	0.00			
Time:	05:56:57	Log-Likelihood:	-1.0374e+05			
No. Observations:	15290	AIC:	2.075e+05			
Df Residuals:	15284	BIC:	2.075e+05			
Df Model:	6					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Number_of_Vehicles	109.7161	1.475	74.368	0.000	106.824	112.608
Gender	139.8706	3.422	40.876	0.000	133.163	146.578
Married	167.3355	3.478	48.113	0.000	160.518	174.153
Fuel_Type	-195.3807	4.983	-39.211	0.000	-205.148	-185.614
YODE_Bucket	-0.9699	0.115	-8.399	0.000	-1.196	-0.744
V_Age_Bucket	10.5786	0.450	23.486	0.000	9.696	11.461
Omnibus:	197.836	Durbin-Watson:	1.269			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	245.423			
Skew:	0.205	Prob(JB):	5.10e-54			
Kurtosis:	3.466	Cond. No.	90.3			

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

OLS with Backward Elimination on X_YODE where Y is dependent variable containing Losses_Capped.

- $R^2 = 0.772$
- $AIC = 2.075$
- Durbin-Watson = 1.269
- “Number_of_Vehicles” is removed by Backward Elimination because of its high p-value.

■ →→→ We selected Model-1 because of its high R^2 value & low AIC value.

CODE:

```
1. import pandas as pd
2.
3. #####
4.
5. # 1. Data Consolidation:
6. dataset = pd.read_csv('Historical_Losses_Data.csv')
7.
8. # 1.1 Gathering Information about the dataset
9. dataset.shape
10.     dataset.head(6)
11.     dataset.tail(6)
12.     dataset.info()
13.
14. # 1.2 Removing NON-Significant Variables
15.     dataset = dataset.drop(['Policy_Number'], axis=1)
16.
17. # 1.3 Creating a list to gather variable/feature name
18.     columns = dataset.columns.tolist()
19.
20. #####
21.
22. # 2. Exploratory Data Analysis -EDA Report
23.     import seaborn as sns
24.     import matplotlib.pyplot as plt
25.     import pandasql as ps
26.     plt.style.use('seaborn-whitegrid')  ### setting up Grid-Lines Style
27.
28. ##### 2.1 Univariate Analysis: #####
29.
30. # 2.1.1. finding ranges and catgories in variables
31.     dataset.Age.sort_values().unique()
32.     dataset.Years_of_Driving_Experience.sort_values().unique()
33.     dataset.Number_of_Vehicles.sort_values().unique()
34.     dataset.Gender.sort_values().unique()
35.     dataset.Married.sort_values().unique()
36.     dataset.Vehicle_Age.sort_values().unique()
37.     dataset.Fuel_Type.sort_values().unique()
38.     dataset.Losses.sort_values().unique()
39.
40. # 2.1.2. Description:
41.     Description = dataset.describe()
42.     skew=dataset.skew()
43.
44. # 2.1.3. Histogram of dataset
45.     dataset.hist(bins=500, figsize=(20,15))
46.     plt.savefig('Historical_Data(Histogram).png')
47.     plt.show()
```



```

48.
49.     # taking a seperate look at 'Losses' variable because of high positive skewness:
50.     dataset['Losses'].hist(bins=500, figsize=(10,7))
51.     plt.title('Losses (Histogram)\n')
52.     plt.savefig('Losses(Histogram).png')
53.     plt.show()
54.
55.     # upper-capping 'losses' at 1200 according to industrial standards to deal with
skewness:
56.     Losses_Capped = dataset['Losses'].clip(upper=1200)
57.     Losses_Capped.hist(bins=500, figsize=(10,7))
58.     plt.title('Losses Capped (Histogram)\n')
59.     plt.savefig('Losses_Capped(Histogram).png')
60.     plt.show()
61.     Losses_Capped.skew()
62.
63.     # making a single DataFrame of 'dataset' and 'Losses_Capped' for further analysis:
64.     q1 = """ SELECT *,
65.             (CASE WHEN Losses > 1200 THEN 1200 ELSE Losses END)
66.             AS Losses_Capped
67.             FROM dataset """
68.     Capped_Losses_dataset = ps.sqlldf(q1, locals())
69.
70.     # 2.1.4. Pie-Plots:
71.     # "Number_of_Vehicles" variable:
72.     Number_of_Vehicles_Count = pd.DataFrame(
73.         dataset.Number_of_Vehicles.value_counts().reset_index())
74.     Number_of_Vehicles_Count.columns=['Number of Vehicles','Count']
75.     plt.pie(Number_of_Vehicles_Count['Count'],
76.             labels=Number_of_Vehicles_Count['Number of Vehicles'],
77.             autopct='%.1f%%', startangle=90)
78.     plt.title('Number of Vehicles')
79.     plt.savefig('Number_of_Vehicles_Count(Pie-Plot).png')
80.     plt.show()
81.
82.     # taking a seperate view at charachter categorical data (Demographical Variables):
83.
84.     #Gender:
85.     Gender_Count = pd.DataFrame(dataset.Gender.value_counts().reset_index())
86.     Gender_Count.columns=['Gender','Count']
87.     plt.pie(Gender_Count['Count'], labels=Gender_Count['Gender'],
88.             autopct='%.1f%%', startangle=90)
89.     plt.title('Gender Count')
90.     plt.savefig('Gender_Count(Pie-Plot).png')
91.     plt.show()
92.
93.     #Married:
94.     Married_Count = pd.DataFrame(dataset.Married.value_counts().reset_index())
95.     Married_Count.columns=['Married?', 'Count']

```

```

96. plt.pie(Married_Count['Count'], labels=Married_Count['Married?'],
97.         colors=['Yellow','Cyan'], autopct='%.1f%%', startangle=90)
98. plt.title('Married?')
99. plt.savefig('Married_Count(Pie-Plot).png')
100. plt.show()
101.
102. #Fuel_Type:
103. Fuel_Type_Count=pd.DataFrame(dataset.Fuel_Type.value_counts().reset_index())
104. Fuel_Type_Count.columns=['Fuel_Type','Count']
105. plt.pie(Fuel_Type_Count['Count'], labels=Fuel_Type_Count['Fuel_Type'],
106.         colors=['y','Red'], autopct='%.1f%%', startangle=90)
107. plt.title('Fuel Type')
108. plt.savefig('Fuel_Type_Count(Pie-Plot).png')
109. plt.show()
110.
111. # 2.1.5. Box-Plot of the variable for outlier detection:
112. dataset.boxplot(column=['Age','Years_of_Driving_Experience'], figsize=(5,10))
113. plt.savefig('Historical_dataset(Box-Plot)1.png', bbox_inches='tight')
114. plt.show()
115. dataset.boxplot(column=['Number_of_Vehicles','Vehicle_Age'], figsize=(5,10))
116. plt.savefig('Historical_dataset(Box-Plot)2.png', bbox_inches='tight')
117. plt.show()
118.
119.
120. ##### 2.2 Bivariate Analysis: #####
121.
122. # Dependent Variables are Losses & Losses_Capped
123. # although we are performing analysis on Losses_Capped
124.
125. #2.2.0. Correlation & Pair_Plots:
126. sns.heatmap(dataset.corr(), annot=True)
127. plt.title('Hiatorical Data Correlation (Heatmap)\n')
128. plt.savefig('Historical_Data_Correlation(Heatmap).png', bbox_inches='tight')
129. plt.show()
130.
131. sns.pairplot(dataset, diag_kind='kde')
132. plt.savefig('Historical_Data(Pair-Plot).png')
133. plt.show()
134.
135. # 2.2.1 Age vs Losses_Capped (Average) Analysing:
136. q1 = """ SELECT Age,
137.           AVG(Losses_Capped)
138.           AS Losses_Capped_Avg
139.           FROM Capped_Losses_dataset
140.           GROUP BY Age
141.           ORDER BY Age """
142.
143. Age_vs_Losses_Cap = ps.sqlldf(q1, locals())
144. plt.plot(Age_vs_Losses_Cap['Age'], Age_vs_Losses_Cap['Losses_Capped_Avg'], 'g.-')

```

```

145. plt.title('Age -vs- Losses (Capped)[Avg]\n')
146. plt.xlabel('Age')
147. plt.ylabel('Losses_Capped_Avg')
148. plt.savefig('Age-vs-Losses(Plot).png', bbox_inches='tight')
149. plt.show()
150.
151. # 2.2.2 Years_of_Driving_Experience vs Losses_Capped(Avg) Analysing:
152. # YODE is used as an abbreviation of Years_of_Driving_Experience here.
153. q1= """ SELECT Years_of_Driving_Experience,
154.         AVG(Losses_Capped)
155.         AS Losses_Capped_Avg
156.         FROM Capped_Losses_dataset
157.         GROUP BY Years_of_Driving_Experience
158.         ORDER BY Years_of_Driving_Experience """
159.
160. YODE_vs_Losses_Cap = ps.sqldf(q1, locals())
161.
162. plt.plot(YODE_vs_Losses_Cap['Years_of_Driving_Experience'],
163.         YODE_vs_Losses_Cap['Losses_Capped_Avg'],
164.         color='Brown', marker='.',linestyle='-')
165.
166. plt.title('Years of Driving Experience -vs- Losses(Capped)[Avg]\n')
167. plt.xlabel('Years of Driving Experience')
168. plt.ylabel('Losses_Capped_Avg')
169. plt.savefig('Years_of_Driving_Experience-vs-Losses(Plot).png', bbox_inches='tight')
170. plt.show()
171.
172. # 2.2.3 Number_of_Vehicles vs Losses_Capped(Avg) Analysing:
173. q1= """ SELECT Number_of_Vehicles,
174.         AVG(Losses_Capped)
175.         AS Losses_Capped_Avg
176.         FROM Capped_Losses_dataset
177.         GROUP BY Number_of_Vehicles
178.         ORDER BY Number_of_Vehicles """
179.
180. Number_of_Vehicles_vs_Losses_Cap = ps.sqldf(q1, locals())
181.
182. plt.plot(Number_of_Vehicles_vs_Losses_Cap['Number_of_Vehicles'],
183.         Number_of_Vehicles_vs_Losses_Cap['Losses_Capped_Avg'], 'bo-')
184. plt.title('Number of Vehicles -vs- Losses (Capped)[Avg]\n')
185. plt.xlabel('Number of Vehicles')
186. plt.ylabel('Losses_Capped_Avg')
187. plt.savefig('Number_of_Vehicles-vs-Losses(Plot).png', bbox_inches='tight')
188. plt.show()
189.
190. # 2.2.4 Gender vs Losses_Capped(Avg) Analysis:
191. q1= """ SELECT Gender,
192.         AVG(Losses_Capped)
193.         AS Losses_Capped_Avg

```

```

194.         FROM Capped_Losses_dataset
195.         GROUP BY Gender
196.         ORDER BY Gender """"
197.
198. Gender_vs_Losses_Cap = ps.sqldf(q1, locals())
199. plt.plot(Gender_vs_Losses_Cap['Gender'],
200.          Gender_vs_Losses_Cap['Losses_Capped_Avg'], 'co-.')
201. plt.title('Gender -vs- Losses (Capped)[Avg]\n')
202. plt.xlabel('Gender')
203. plt.ylabel('Losses_Capped_Avg')
204. plt.savefig('Gender-vs-Losses(Plot).png', bbox_inches='tight')
205. plt.show()
206.
207. #2.2.5 Married vs Losses_Capped(Avg) Analysis:
208. q1= """" SELECT Married,
209.          AVG(Losses_Capped)
210.          AS Losses_Capped_Avg
211.          FROM Capped_Losses_dataset
212.          GROUP BY Married
213.          ORDER BY Married """"
214.
215. Married_vs_Losses_Cap = ps.sqldf(q1, locals())
216. plt.plot(Married_vs_Losses_Cap['Married'],
217.          Married_vs_Losses_Cap['Losses_Capped_Avg'], 'ro-.')
218. plt.title('Married? -vs- Losses (Capped)[Avg]\n')
219. plt.xlabel('Married?')
220. plt.ylabel('Losses_Capped_Avg')
221. plt.savefig('Married_vs_Losses(Plot).png', bbox_inches='tight')
222. plt.show()
223.
224. # 2.2.6 Vehicle_Age vs Losses_Capped(Avg) Analysis:
225. q1= """" SELECT Vehicle_Age,
226.          AVG(Losses_Capped)
227.          AS Losses_Capped_Avg
228.          FROM Capped_Losses_dataset
229.          GROUP BY Vehicle_Age
230.          ORDER BY Vehicle_Age """"
231.
232. Vehicle_Age_vs_Losses_Cap = ps.sqldf(q1, locals())
233. plt.plot(Vehicle_Age_vs_Losses_Cap['Vehicle_Age'],
234.          Vehicle_Age_vs_Losses_Cap['Losses_Capped_Avg'], 'mo-.')
235. plt.title('Vehicle Age -vs- Losses (capped)[Avg]\n')
236. plt.xlabel('Vehicle Age')
237. plt.ylabel('Losses_Capped_Avg')
238. plt.savefig('Vehicle_Age-vs-Losses(Plot).png', bbox_inches='tight')
239. plt.show()
240.
241. # 2.2.7 Fuel_Type vs Losses_Capped(Avg) Analysis:
242. q1= """" SELECT Fuel_Type,

```

```

243.     AVG(Losses_Capped)
244.     AS Losses_Capped_Avg
245. FROM Capped_Losses_dataset
246. GROUP BY Fuel_Type
247. ORDER BY Fuel_Type """"
248.
249. Fuel_Type_vs_Losses_Cap = ps.sqldf(q1, locals())
250. plt.plot(Fuel_Type_vs_Losses_Cap['Fuel_Type'],
251.          Fuel_Type_vs_Losses_Cap['Losses_Capped_Avg'], 'yo-.')
252. plt.title('Fuel Type -vs- Losses (Capped)[Avg]\n')
253. plt.xlabel('Fuel Type')
254. plt.ylabel('Losses_Capped_Avg')
255. plt.savefig('Fuel_Type-vs-Losses(Plot).png', bbox_inches='tight')
256. plt.show()
257.
258. # EDA Report Ends.
259. #####
260.
261. ### 2.3 EDA Conclusion: ###
262.
263. """" From the above EDA Report we get to know that
264.     “Age”, “Years_of_Driving_Experience” and “Vehicle_Age” needs bucketing""""
265.
266. # 2.4 Bucketing:
267. # 2.4.1. Age Bucketing:
268.
269. q1 = """" SELECT *,
270.     (CASE
271.         WHEN Age BETWEEN 16 AND 25 THEN 21
272.         WHEN Age BETWEEN 26 AND 59 THEN 43
273.         WHEN Age BETWEEN 60 AND 70 THEN 65
274.         ELSE Age
275.     END)
276.     AS Age_Bucket
277. FROM Capped_Losses_dataset
278. ORDER BY Age""""
279.
280. CLD_Age_Bucket = ps.sqldf(q1, locals())
281.
282. #2.4.2. Years_of_Driving_Experience Bucketing:
283.
284. q1 = """" SELECT *,
285.     (CASE
286.         WHEN Years_of_Driving_Experience BETWEEN 0 AND 8 THEN 4
287.         WHEN Years_of_Driving_Experience BETWEEN 9 AND 40 THEN 25
288.         WHEN Years_of_Driving_Experience BETWEEN 41 and 53 THEN 47
289.         ELSE Years_of_Driving_Experience
290.     END)
291.     AS YODE_Bucket

```

```

292.         FROM CLD_Age_Bucket
293.         ORDER BY Years_of_Driving_Experience""
294.
295. CLD_YODE_Bucket = ps.sqldf(q1, locals())
296.
297. #2.4.3. Vehicle_Age Bucketing:
298.
299. q1= "" SELECT *,
300.     (CASE
301.         WHEN Vehicle_Age BETWEEN 0 AND 5 THEN 3
302.         WHEN Vehicle_Age BETWEEN 6 AND 10 THEN 8
303.         WHEN Vehicle_Age BETWEEN 11 AND 15 THEN 13
304.         ELSE Vehicle_Age
305.     END)
306.     AS V_Age_Bucket
307.     FROM CLD_YODE_Bucket
308.     ORDER BY Age""
309.
310. CLD_V_Age_Bucket = ps.sqldf(q1, locals())
311.
312.
313. # dropping variables that we have bucketed & capped
314. dataset_final = CLD_V_Age_Bucket.drop(columns=['Age',
315.         'Years_of_Driving_Experience',
316.         'Vehicle_Age','Losses'], axis=1)
317.
318. #####
    ###
319.
320. ##### 3 Confirmatory Data Analysis: #####
321.
322. # 3.1 gathering information about 'dataset_final':
323. Description_final = dataset_final.describe()
324. Skew_final = dataset_final.skew()
325.
326. # 3.2 Bivariate Analysis (CDA):
327.
328. # 3.2.1 Correlation Check:
329. sns.heatmap(dataset_final.corr(), annot=True)
330. plt.title('Final Dataset Correlation\n')
331. plt.savefig('dataset_final_correlation(Heatmap).png', bbox_inches='tight')
332. plt.show()
333.
334. # "Years_of_Driving_Experience" & "Age" still have high correlation even after bucketing
335. # there should be two models in model preperation because of this correlation
336.
337. # 3.2.2 Age_Bucket vs Losses_Capped(Avg) :
338.

```

```

339. q1 = """ SELECT Age_Bucket,
340.         AVG(Losses_Capped)
341.         AS Losses_Capped_Avg
342.         FROM dataset_final
343.         GROUP BY Age_Bucket
344.         ORDER BY Age_Bucket """
345.
346. Age_Bucket_vs_Losses_Cap = ps.sqlldf(q1, locals())
347. plt.plot(Age_Bucket_vs_Losses_Cap['Age_Bucket'],
348.         Age_Bucket_vs_Losses_Cap['Losses_Capped_Avg'], 'g.-')
349. plt.title('Age Bucket -vs- Losses (Capped)[Avg]\n')
350. plt.xlabel('Age (Bucket)')
351. plt.ylabel('Losses_Capped_Avg')
352. plt.savefig('Age_Bucket-vs-Losses_Cap(Plot).png', bbox_inches='tight')
353. plt.show()
354.
355. #3.2.3 YODE_Bucket vs Losses_Capped(Avg) :
356.
357. q1 = """ SELECT YODE_Bucket,
358.         AVG(Losses_Capped)
359.         AS Losses_Capped_Avg
360.         FROM dataset_final
361.         GROUP BY YODE_bucket
362.         ORDER BY YODE_Bucket """
363.
364. YODE_Bucket_vs_Losses_Cap = ps.sqlldf(q1, locals())
365. plt.plot(YODE_Bucket_vs_Losses_Cap['YODE_Bucket'],
366.         YODE_Bucket_vs_Losses_Cap['Losses_Capped_Avg'],
367.         color='Brown', marker='.', linestyle='-')
368.
369. plt.title('YODE_Bucket -vs- Losses(Capped)[Avg]\n')
370. plt.xlabel('Years of Driving Experience (Bucket)')
371. plt.ylabel('Losses_Capped_Avg')
372. plt.savefig('YODE_Bucket-vs-Losses_Cap(Plot).png', bbox_inches='tight')
373. plt.show()
374.
375. # 3.2.4 V_Age_Bucket vs Losses_Capped(Avg):
376.
377. q1 = """ SELECT V_Age_Bucket,
378.         AVG(Losses_Capped)
379.         AS Losses_Capped_Avg
380.         FROM dataset_final
381.         GROUP BY V_Age_Bucket
382.         ORDER BY V_Age_Bucket """
383.
384. V_Age_Bucket_vs_Losses_Cap = ps.sqlldf(q1, locals())
385. plt.plot(V_Age_Bucket_vs_Losses_Cap['V_Age_Bucket'],
386.         V_Age_Bucket_vs_Losses_Cap['Losses_Capped_Avg'],
387.         'mo-.')

```

```
388. plt.title('Vehicle Age -vs- Losses (capped)[Avg]\n')
389. plt.xlabel('Vehicle Age (Bucket)')
390. plt.ylabel('Losses_Capped_Avg')
391. plt.savefig('V_Age_Bucket-vs-Losses_Cap(Plot).png', bbox_inches='tight')
392. plt.show()
393.
394. #####
395.
396. #### 4. Data Pre-Processing: #### (on 'dataset_final')
397.
398. dataset_final_OLS = dataset_final.copy() # making a copy
399.
400. # 4.1 Missing Value Treatment:
401. " There are no missing value so skipping this step"
402.
403. # 4.2 Oulier Treatment :
404. " Outlier Values have been treated above with the help of 'Capping' and 'Bucketing' "
405.
406. # 4.3 Variable Transformation :
407. from sklearn.preprocessing import LabelEncoder
408.
409. # 4.3.1 converting categorical data into binary/numerical form:
410. myencoder = LabelEncoder()
411.
412. dataset_final_OLS['Gender'] = myencoder.fit_transform(dataset_final_OLS['Gender'])
413. dataset_final_OLS['Married'] = myencoder.fit_transform(dataset_final_OLS['Married'])
414. dataset_final_OLS['Fuel_Type'] =
    myencoder.fit_transform(dataset_final_OLS['Fuel_Type'])
415.
416. # 4.4 Variable Creation :
417.
418. # Segregation between Independent and Dependent Vaariables:
419. # X = Independent Variables ; Y = Dependent Variables
420.
421. X = pd.DataFrame(dataset_final_OLS.drop(columns='Losses_Capped')).copy()
422. Y = pd.DataFrame(dataset_final_OLS['Losses_Capped']).copy()
423.
424. # from CDA we know "YODE_Bucket" & "Age_Bucket" have correlation nearly to 1
425. # making two seperate models for these variables
426. X_Age = pd.DataFrame(X.drop(columns='YODE_Bucket')).copy()
427. X_YODE = pd.DataFrame(X.drop(columns='Age_Bucket')).copy()
428.
429. #####
430.
431. ##### 5. Model Development: #####
432. import statsmodels.formula.api as smfa
433. import statsmodels.tools as smt
434.
435. # 5.1 Creating Linear Formula:
```



```

436. #  $b_0X_0 + b_1X_1 + b_2X_2 + \dots + b_8X_8$ 
437. # where  $X_0$  should be constant
438. X = smt.add_constant(X)
439.
440. # 5.2 Creating Different Models with OLS (Ordinary Least Square)
441. ##### & using Backward Elimination Approach
442.
443. import selectionprocess as sp ## user generated module
444. # setting a significance level of 5% (i.e, 0.05)
445. sig_level = 0.05
446.
447. # 5.2.1 Model-1 : on X_Age
448. X_Age_Model = Age_Model_Summary =
    sp.BackwardElimination_OLS_DataFrame(X_Age, Y, sig_level)
449.
450. # "Number of_Vehicles" is removed by backward elimination through OLS
451. # because its p-value > significance level
452.
453. font_dict = {'family' : 'monospace',
454.              'size' : 'large',
455.              'weight' : 'semibold'}
456.
457. plt.text(0, 0, str(Age_Model_Summary),font_dict)
458. plt.axis('off')
459. plt.savefig('OLS_report_X_Age_Model.png', bbox_inches='tight')
460. plt.show()
461.
462. # 5.2.2 Model-2 : on X_YODE
463. X_YODE_Model = YODE_Model_Summary =
    sp.BackwardElimination_OLS_DataFrame(X_YODE, Y, sig_level)
464.
465. # "Number of_Vehicles" is removed by backward elimination through OLS
466. # because its p-value > significance level
467.
468. plt.text(0, 0, str(YODE_Model_Summary),font_dict)
469. plt.axis('off')
470. plt.savefig('OLS_report_X_YODE_Model.png', bbox_inches='tight')
471. plt.show()
472.
473.
474. """ selecting Model-1 for further prediction
475. as it has higher  $R^2$  value and lower AIC value """
476.
477. coeff = smfa.OLS(Y, X_Age_Model).fit().params
478.
479. # 5.2.3 Splitting into Train Test Values:
480. from sklearn.model_selection import train_test_split
481.
482. X_Train, X_Test, Y_Train, Y_Test = train_test_split(X_Age_Model, Y,

```

```
483.                                     test_size=0.2, random_state=66)
484.
485.     from sklearn.linear_model import LinearRegression
486.     regressor = LinearRegression()
487.     regressor.fit(X_Train, Y_Train)
```

#####

```
488.     ##### 6. Model Performance: #####
489.
490.     Y_Predicted = regressor.predict(X_Test)
491.
492.     # 6.1 Creating a single Analysis DataFrame:
493.
494.     Y_Predicted_temp = pd.DataFrame(Y_Predicted, columns=['Losses_Predicted'])
495.     Y_Test_temp = pd.DataFrame(Y_Test.reset_index().drop('index',axis=1))
496.
497.     temp = dataset_final.copy()
498.     tempTrain, tempTest=train_test_split(temp,test_size=0.2,random_state=66)
499.
500.     temp = tempTest.drop(columns=['Losses_Capped','YODE_Bucket'])
501.     temp = pd.DataFrame(temp.reset_index().drop('index',axis=1))
502.
503.     Analysis = pd.concat([temp, Y_Test_temp, Y_Predicted_temp], axis=1)
504.     Error_Values = pd.DataFrame(Analysis['Losses_Predicted']-Analysis['Losses_Capped'],
505.                                columns=['Error (Losses_Predicted-Losses_Capped)'])
506.
507.     Analysis = pd.concat([Analysis, Error_Values], axis=1)
508.
509.     #END.
510.
```

1. # Slection Process all in one module

2.

3. import numpy as np

4. import pandas as pd

5. import statsmodels.formula.api as smfa

6.

7. # 1 Backward Elimination: for array-like

8. # independent_array == array-like, independent dataset/variables (normally denoted by X in program)

9. # dependent_array == array-like(1d), dependent dataset/variable

10. # significance level == float, if value of significance is 5% set this value as 0.05

11.

12. # returns(array-like) == moduled or processed final array after removing all

13. # in-significant variables through OLS

14. # returns (summary) == summary of final model as

'statsmodels.iolib.summary.Summary'

15. # also prints final summary

16.

17. def Backward_Elimination_OLS_Array(independent_array, dependent_array, significance):

18. count = len(independent_array[0])

19. for num1 in range(count):

20. OLS_Regressor = smfa.OLS(dependent_array, independent_array).fit()

21. Max_P_value = max(OLS_Regressor.pvalues).astype(float)

22.

23. if Max_P_value > significance:

24. for num2 in range(count-num1):

25. if(OLS_Regressor.pvalues[num2].astype(float) == Max_P_value):

26. independent_array = np.delete(arr=independent_array,

27. obj=num2, axis=1)

28. print(OLS_Regressor.summary())

29. return independent_array, OLS_Regressor.summary()

30.

31.

32. # 2 Backward Elimination: for DataFrame-like

33. # independent_DF == DataFrame-like, independent dataset/variables (normally denoted by X in program)

34. # dependent_DF == DataFrame-like(1d), dependent dataset/variable

35. # significance level == float, if value of significance is 5% set this value as 0.05

36.

37. # returns (DataFrame-like)== moduled or processed final DataFrame after removing all

38. # in-significant variables through OLS

39. # returns (summary) == summary of final model as

'statsmodels.iolib.summary.Summary'

40.

```

41.     # also prints final summary
42.
43.
44.     def BackwardElimination_OLS_DataFrame(independent_DF, dependent_DF,
        significance):
45.         count = len(independent_DF.columns)
46.         for num1 in range(count):
47.             OLS_Regressor = smfa.OLS(dependent_DF, independent_DF).fit()
48.             Max_P_value = max(OLS_Regressor.pvalues)
49.
50.             if Max_P_value > significance:
51.                 for num2 in range(count-num1):
52.                     if(OLS_Regressor.pvalues[num2] == Max_P_value):
53.                         independent_DF = pd.DataFrame.drop(independent_DF,
54.                             columns=independent_DF.columns[num2],
55.                             axis=1)
56.         print(OLS_Regressor.summary())
57.         return independent_DF, OLS_Regressor.summary()
58.

```