

Insurance Claim Status Prediction

Insurance companies take risks over customers. Risk management is a very important aspect of the insurance industry. Insurers consider every quantifiable factor to develop profiles of high and low insurance risks. Insurers collect vast amounts of information about policyholders and analyse the data. As a Data scientist in an insurance company, you need to analyse the available data and predict whether to approve the insurance or not.

Prepared By Humera Shaikh

Importing Library

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import skew
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.metrics import roc_auc_score, roc_curve
from sklearn.feature_selection import chi2
from sklearn.feature_selection import f_classif
from sklearn.feature_selection import SelectKBest
from sklearn.ensemble import VotingClassifier

from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.svm import SVC
from sklearn.svm import LinearSVC
import warnings
warnings.filterwarnings("ignore")

In [2]: df=pd.read_csv("insurane claim.csv")

In [3]: df.head()
```

	ID	Agency	Agency Type	Distribution Channel	Product Name	Claim	Duration	Destination	Net Sales	Commission (in value)	Gender	Age
0	3433	CWT	Travel Agency	Online	Rental Vehicle Excess Insurance	0	7	MALAYSIA	0.0	17.82	NaN	31
1	4339	EPX	Travel Agency	Online	Cancellation Plan	0	85	SINGAPORE	69.0	0.00	NaN	36
2	34590	CWT	Travel Agency	Online	Rental Vehicle Excess Insurance	0	11	MALAYSIA	19.8	11.88	NaN	75
3	55816	EPX	Travel Agency	Online	2 way Comprehensive Plan	0	16	INDONESIA	20.0	0.00	NaN	32
4	13816	EPX	Travel Agency	Online	Cancellation Plan	0	10	KOREA, REPUBLIC OF	15.0	0.00	NaN	29

```
In [4]: df.shape
Out[4]: (50553, 12)

EDA

In [5]: # Checking the weightage of null values.
(df['Gender'].isnull().sum()/(df.shape[0]*100)).round(2)
```

Out[5]: 71.12

We can see 71% of the data does not have Gender updated, there are no other features that will help us find out whether the individual is male or female. Thus we will drop this feature

We can see from data ID column is not useful data,hnce we will drop ID column

```
In [6]: df.drop(["Gender"],axis=1,inplace=True)

In [7]: df.drop(["ID"],axis=1,inplace=True)
```

```
In [8]: df.describe()
```

	Claim	Duration	Net Sales	Commission (in value)	Age
count	50553.000000	50553.000000	50553.000000	50553.000000	50553.000000
mean	0.014658	49.425969	40.800977	9.83809	40.011236
std	0.120180	101.434647	48.899683	19.91004	14.076566
min	0.000000	-2.000000	-389.000000	0.000000	0.000000
25%	0.000000	9.000000	18.000000	0.000000	35.000000
50%	0.000000	22.000000	26.500000	0.000000	36.000000
75%	0.000000	53.000000	48.000000	11.500000	44.000000
max	1.000000	4881.000000	810.000000	283.500000	118.000000

Duration is negative so checking destination for replace negative value

```
In [9]: df[(df["Duration"]<0) & (df["Destination"])]
```

	Agency	Agency Type	Distribution Channel	Product Name	Claim	Duration	Destination	Net Sales	Commission (in value)	Age
4063	IJZI	Airlines	Online	Basic Plan	0	-1	BRUNEI DARUSSALAM	18.0	6.3	118
38935	IJZI	Airlines	Online	Basic Plan	0	-1	INDONESIA	18.0	6.3	118
48367	IJZI	Airlines	Online	Basic Plan	0	-2	BANGLADESH	22.0	7.7	118

```
In [10]: #BANGLADESH, INDONESIA, BRUNEI DARUSSALAM has negative values
#replacing this negative value with mean value of its own
```

```
In [11]: ban=df[(df["Duration"] & (df["Destination"]=="BANGLADESH"))]
ban["Duration"].mean()
```

Out[11]: 15.958133333333334

```
In [12]: ind=df[(df["Duration"]) & (df["Destination"]=="INDONESIA")]
ind["Duration"].mean()
```

Out[12]: 31.3520336058906

```
In [13]: BRUNEI=df[(df["Duration"]) & (df["Destination"]=="BRUNEI DARUSSALAM")]
BRUNEI["Duration"].mean()
```

Out[13]: 22.37012987012987

```
In [14]: des=df[(df["Duration"]==0) & (df["Destination"])]
des
```

	Agency	Agency Type	Distribution Channel	Product Name	Claim	Duration	Destination	Net Sales	Commission (in value)	Age
181	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
314	RAB	Airlines	Online	Value Plan	0	0	BRUNEI DARUSSALAM	15.00	6.00	53
1864	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
3068	C2B	Airlines	Online	Silver Plan	0	0	SINGAPORE	0.00	5.63	51
4282	LWC	Travel Agency	Online	Single Trip Travel Protect Gold	0	0	MALAYSIA	27.00	17.55	56
7324	RAB	Airlines	Online	Value Plan	0	0	BRUNEI DARUSSALAM	15.00	6.00	24
7482	IJZI	Airlines	Online	Basic Plan	0	0	THAILAND	18.00	6.30	118
7865	IJZI	Airlines	Online	Basic Plan	0	0	THAILAND	18.00	6.30	118
8205	JWT	Airlines	Online	Value Plan	0	0	INDIA	62.00	24.80	118
8512	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
8897	JWT	Airlines	Online	Value Plan	0	0	INDIA	62.00	24.80	118
9501	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
10465	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
10558	JWT	Airlines	Online	Value Plan	0	0	INDIA	62.00	24.80	118
11515	IJZI	Airlines	Online	Basic Plan	0	0	THAILAND	18.00	6.30	118
12579	CSB	Airlines	Online	Ticket Protector	0	0	SINGAPORE	5.80	1.63	48
12924	C2B	Airlines	Online	Bronze Plan	0	0	SINGAPORE	-14.40	3.60	49
16251	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	1.80	0.50	48
17641	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	0.51	0.14	48
17741	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	1.35	0.38	48
18198	JWT	Airlines	Online	Value Plan	0	0	INDIA	62.00	24.80	118
18299	RAB	Airlines	Online	Value Plan	0	0	BRUNEI DARUSSALAM	15.00	6.00	33
19419	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	14.40	4.04	48
20538	JWT	Airlines	Online	Value Plan	0	0	INDIA	62.00	24.80	118
21442	C2B	Airlines	Online	Bronze Plan	0	0	SINGAPORE	27.00	6.75	23
23561	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	1.93	0.54	48
23606	IJZI	Airlines	Online	Basic Plan	0	0	MYANMAR	18.00	6.30	24
23766	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	2.44	0.68	48
24505	JWT	Airlines	Online	Basic Plan	0	0	MYANMAR	18.00	6.30	26
25266	CWT	Travel Agency	Online	Rental Vehicle Excess Insurance	0	0	MALAYSIA	9.90	5.94	33
25785	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	4.44	1.25	48
26076	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	0.84	0.24	48
27002	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
30597	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	3.99	1.12	48
31258	EPX	Travel Agency	Offline	1 way Comprehensive Plan	0	0	MYANMAR	20.00	0.00	33
31932	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	2.86	0.80	48
33674	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	4.88	1.37	48
34178	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
34511	JWT	Airlines	Online	Value Plan	0	0	CHINA	22.00	7.70	34
36144	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
36403	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
39014	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
39027	IJZI	Airlines	Online	Basic Plan	0	0	INDONESIA	18.00	6.30	118
39318	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
41302	JWT	Airlines	Online	Value Plan	0	0	INDIA	62.00	24.80	118
41861	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	37.20	118
43410	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	9.77	2.74	48
43464	RAB	Airlines	Online	Value Plan	0	0	BRUNEI DARUSSALAM	15.00	6.00	24
43485	IJZI	Airlines	Online	Basic Plan	0	0	VIET NAM	18.00	6.30	58
43719	SSI	Airlines	Online	Ticket Protector	0	0	SINGAPORE	1.03	0.29	48
45395	CWT	Travel Agency	Online	Rental Vehicle Excess Insurance	0	0	UNITED STATES	9.90	5.94	59
45474	RAB	Airlines	Online	Value Plan	0	0	BRUNEI DARUSSALAM	15.00	6.00	24
47625	LWC	Travel Agency	Online	Single Trip Travel Protect Gold	0	0	INDIA	35.25	22.91	34
50357	C2B	Airlines	Online	Bronze Plan	0	0	SINGAPORE	14.40	3.60	49

```
In [15]: des["Destination"].mode()
Out[15]: 0 INDIA
dtype: object
```

```
In [16]: ind=df[(df["Duration"] & (df["Destination"]=="INDIA*"))]
ind["Duration"].mean()
```

Out[16]: 30.41387024608501

```
In [17]: #duration has negative value
for i in range(len(df)):
    if df["Duration"][i]<0:
        print(df["Duration"].iloc[i])
```

```
-1
-1
-2

In [18]: df[(df["Age"] >100)]
```

	Agency	Agency Type	Distribution Channel	Product Name	Claim	Duration	Destination	Net Sales	Commission (in value)	Age
90	JWT	Airlines	Online	Value Plan	0	58	INDIA	78.0	31.20	118
108	JWT	Airlines	Online	Value Plan	0	15	INDIA	31.0	12.40	118
140	JWT	Airlines	Online	Value Plan	0	8	INDIA	39.0	15.60	118
153	JWT	Airlines	Online	Value Plan	0	4	INDIA	78.0	31.20	118
181	JWT	Airlines	Online	Value Plan	0	0	INDIA	31.00	12.40	118
...	
50158	JWT	Airlines	Online	Value Plan	0	41	INDIA	60.0	24.00	118
50179	JWT	Airlines	Online	Value Plan	0	62	INDIA	31.0	12.40	118
50230	JWT	Airlines	Online	Value Plan	0	15	INDIA	31.0	12.40	118
50429	IJZI	Airlines	Online	Basic Plan	0	19	SRI LANKA	35.0	12.25	118
50478	CCR	Travel Agency	Offline	Comprehensive Plan	0	6	THAILAND	29.0	9.57	118

795 rows × 10 columns

```
In [19]: df["Age"].mean()
Out[19]: 40.011235732795285
```

#We can see from data there is 795 rows so we replace age with mean value 40.011235732795285

```
In [20]: #we can see from data there is 795 rows so we replace age with mean value 40.011235732795285
df.loc[(df["Age"] > 100, "Age") = 40
```

```
In [21]: df[(df["Age"] >100)]
```

	Agency	Agency Type	Distribution Channel	Product Name	Claim	Duration	Destination	Net Sales	Commission (in value)	Age
48367	IJZI	Airlines	Online	Basic Plan	0	-2	BANGLADESH	22.0	7.7	40

```
In [23]: # we replace duration -1 value with 27 INDONESIA,BRUNEI DARUSSALAM travel duration mean time is 27
# we replace duration -2 value with 16 Bangladesh duration mean time is 16
```

```
df.loc[(df["Duration"]!=-1, "Duration") #27
df.loc[(df["Duration"]!=-2, "Duration") #16
```

```
df[(df["Duration"]!=-1)]
```

	Agency	Agency Type	Distribution Channel	Product Name	Claim	Duration	Destination	Net Sales	Commission (in value)	Age
48367	IJZI	Airlines	Online	Basic Plan	0	-2	BANGLADESH	22.0	7.7	40

```
In [25]: # we replace duration 0 value with 30 because mode of 0 value destination is india,and indias mean is 30
df[(df["Duration"]==0, "Duration") #30
```

```
In [26]: df[(df["Duration"]==0)]
```

	Agency	Agency Type	Distribution Channel	Product Name	Claim	Duration	Destination	Net Sales	Commission (in value)	Age
0	CWT	Travel Agency	Online	Rental Vehicle Excess Insurance	0	7	MALAYSIA	0.0	17.82	31
1	EPX	Travel Agency	Online	Cancellation Plan	0	85	SINGAPORE	69.0	0.00	36
2	CWT	Travel Agency	Online	Rental Vehicle Excess Insurance	0	11	MALAYSIA	19.8	11.88	75
3	EPX	Travel Agency	Online	2 way Comprehensive Plan	0	16	INDONESIA	20.0	0.00	32
4	EPX	Travel Agency	Online	Cancellation Plan	0	10	KOREA, REPUBLIC OF	15.0	0.00	29

```
In [28]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50553 entries, 0 to 50552
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Agency                50553 non-null object
1   Agency Type           50553 non-null object
2   Distribution Channel   50553 non-null object
3   Product Name          50553 non-null object
4   Claim                 50553 non-null int64
5   Duration              50553 non-null int64
6   Destination           50553 non-null object
7   Net Sales             50553 non-null float64
8   Commission (in value) 50553 non-null float64
9   Age                  50553 non-null int64
dtypes: float64(2), int64(3), object(5)
memory usage: 3.9+ MB
```

```
In [30]: #changing date of claim int to objctive
df["Claim"]=df["Claim"].astype('int')
```

```
In [30]: #net sale wont be negative so removing those record
df=df[(df["Net Sales"] > 0)]
df["Net Sales"].describe().T
```

count	48536.000000
mean	43.148207
std	47.939269
min	0.070000
25%	18.800000
50%	28.000000
75%	49.500000
max	810.000000
Name: Net Sales, dtype: float64	

```
In [31]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 48536 entries, 1 to 50552
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Agency                48536 non-null object
1   Agency Type           48536 non-null object
2   Distribution Channel   48536 non-null object
3   Product Name          48536 non-null object
4   Claim                 48536 non-null int32
5   Duration              48536 non-null int64
6   Destination           48536 non-null object
7   Net Sales             48536 non-null float64
8   Commission (in value) 48536 non-null float64
9   Age                  48536 non-null int64
dtypes: float64(2), int32(1), int64(2), object(5)
memory usage: 3.9+ MB
```

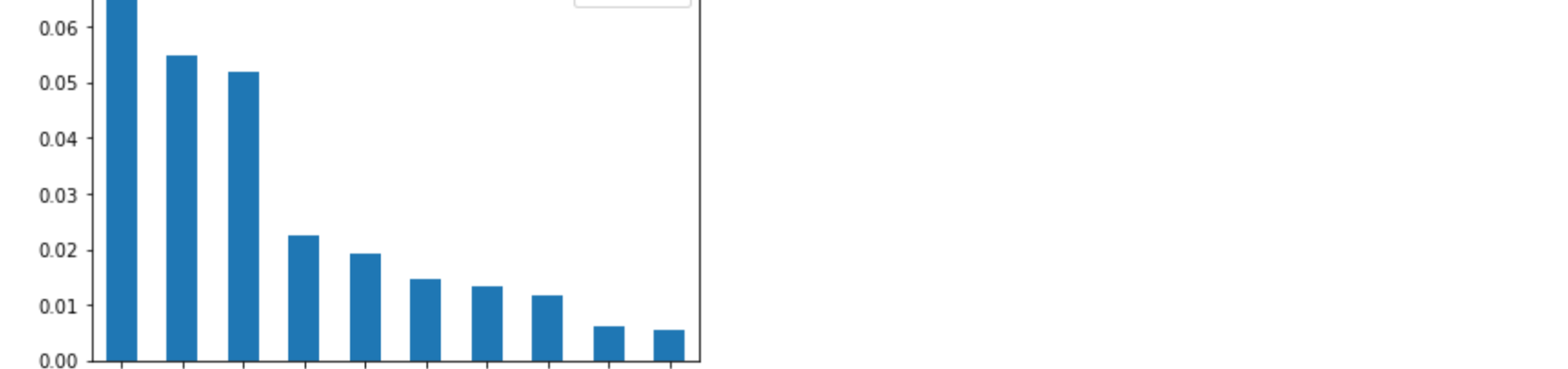
Data visualization

Agencies with maximum claims

```
In [32]: #Plotting diffrant diffrant graphs for understanding
# Finding the Agencies with the Maximum Claims
Claim_mean = (df.groupby('Agency', as_index=False)['Claim'].mean()).nlargest(10,'Claim')
Claim_mean
```

	Agency	Claim
2	C2B	0.067726
11	LWC	0.055028
15	TTW	0.051948
10	KML	0.022364
4	CCR	0.019555
5	CSR	0.014706
3	CBH	0.013333
6	CWT	0.011810
13	SSI	0.006010
7	EPX	0.005468

```
In [33]: ax = Claim_mean.plot.bar(x='Agency', y='Claim', rot=0)
```



Out(50):

	Agency	Agency Type	Distribution Channel	Product Name	Destination
	1	7	1	1	10 79
	2	6	1	1	16 56
	3	7	1	1	1 38
	4	7	1	1	10 47
	5	6	1	1	16 88
...
50548	6	1	1	1	16 65
50549	7	1	1	1	0 38
50550	10	1	1	24	38
50551	7	1	1	10	7
50552	7	1	1	0	79

48536 rows × 5 columns

In (51):

df_newmpd.concat([df_num,df_cat],axis=1)

In (52):

df_new

Out(52):

	Duration	Net Sales	Commission (in value)	Age	Claim	Agency	Agency Type	Distribution Channel	Product Name	Destination
1	9.219544	8.306624	0.000000	6.000000	0	7	1	1	10	79
2	3.316625	4.449719	3.446738	8.660254	0	6	1	1	16	56
3	4.000000	4.472136	0.000000	5.656854	0	7	1	1	1	38
4	3.162278	3.872983	0.000000	5.385165	0	7	1	1	10	47
5	8.000000	7.035624	5.449771	6.000000	0	6	1	1	16	88
...

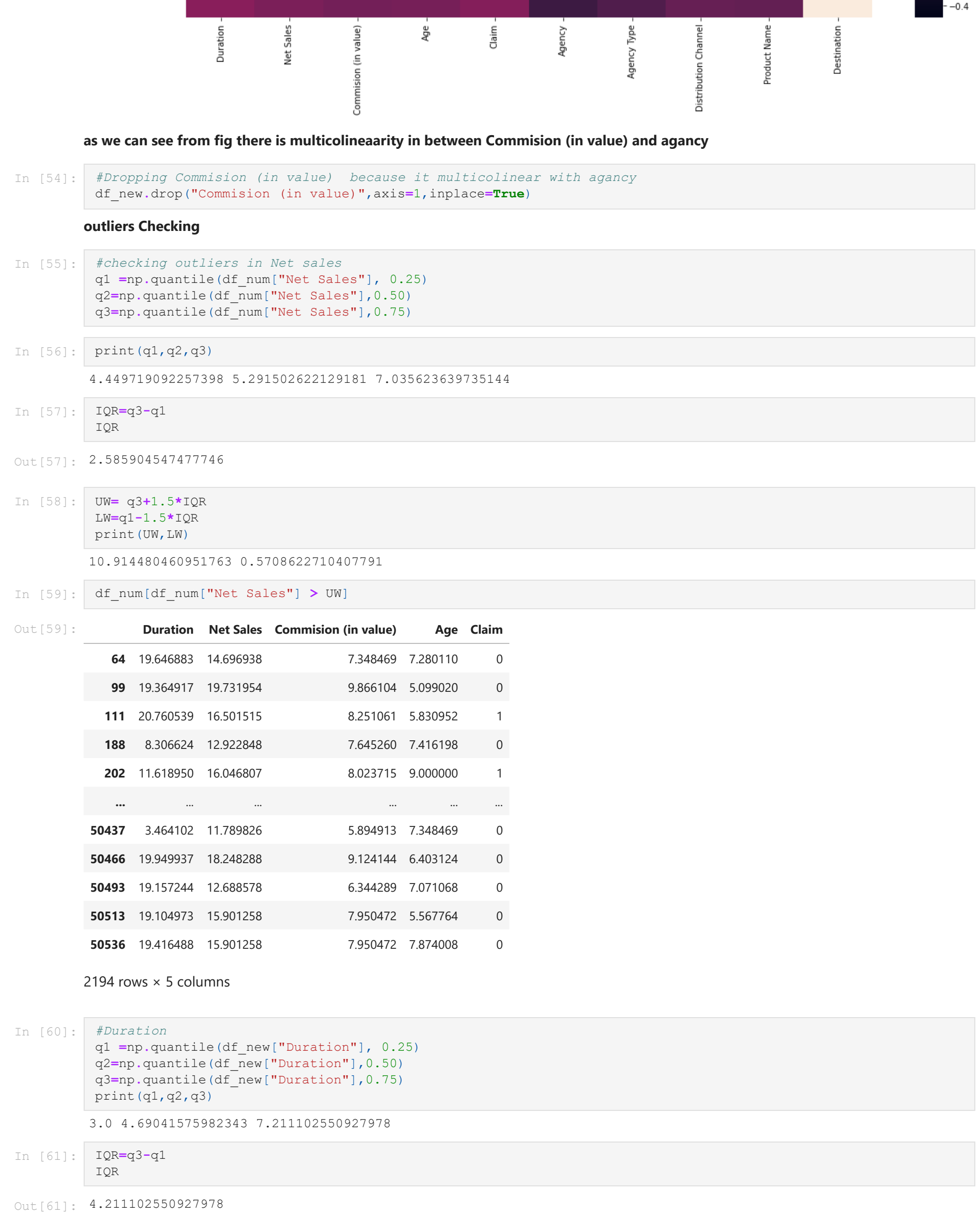
48536 rows x 5 columns

df_new=pd.concat((df_num,df_cat),axis=1)										
df_new										
Duration Net Sales Commission (in value) Age Claim Agency Agency Type Distribution Channel Product Name Destination										
1	9.219544	8.306624	0.000000	6.000000	0	7	1	1	10	79
2	3.316625	4.449719	3.446738	6.660254	0	6	1	1	16	56

3	4.000000	4.472136	0.000000	5.656854	0	7	1	1	1	38
4	3.162278	3.872983	0.000000	5.385165	0	7	1	1	1	47
5	8.000000	7.035624	5.449771	6.000000	0	6	1	1	16	88
...										
50548	3.316625	9.949874	7.707140	7.071068	0	7	1	1	16	65
50549	6.928203	4.472136	0.000000	6.000000	0	6	1	1	1	38
50550	4.898979	5.291503	3.261901	6.928203	0	10	1	1	24	38
50551	5.000000	3.872983	0.000000	6.000000	0	7	1	1	1	7
50552	3.000000	4.472136	0.000000	7.681146	0	7	1	1	0	79

48536 rows x 10 columns

#plotting heatmap										
plt.show()										
plt.figure(figsize=(30,12))										
sns.heatmap(df_new.corr(),annot=True)										



as we can see from fig there is multicollinearity in between Commission (in value) and agency

#Dropping Commission (in value) because it multicollinear with agency										
df_new.drop("Commission (in value)",axis=1,inplace=True)										

outliers Checking

#checking outliers in Net sales										
q1=np.quantile(df_new["Net Sales"], 0.25)										
q2=np.quantile(df_new["Net Sales"], 0.50)										
q3=np.quantile(df_new["Net Sales"], 0.75)										

print(q1,q2,q3)

4.444719092257398 5.291502622129181 7.035623639735144

IQR=q3-q1										
IQR										

2.585904547477746

IQR=q1-1.5*IQR										
IQRq1=1.5*IQR										
print(IQR,IQRq1)										

10.914490460951763 0.5708622710407791

50525	16.309506	6.292853	6.405124	0	6	1	1	16	36
50536	19.416488	15.901258	7.874008	0	2	0	1	4	79
50546	13.784049	3.162278	6.000000	0	7	1	1	10	88

2502 rows x 9 columns

```
In [64]: df_new[df_new["Duration"] < 110]
```

188	8.306624	12.922848	7.645260	7.416198	0						
202	16.18950	16.046807	8.023715	9.000000	1						

2	3.316625	4.449719	8.660254	0	6	1	1	16	56
3	4.000000	4.472136	5.656854	0	7	1	1	1	38
4	3.316625	4.449719	8.660254	0	6	1	1	16	56

50466	19.949937	18.248828	9.124144	6.403124	0						
50493	19.157244	12.688578	6.344289	7.071068	0						

50513	19.104973	15.901258	7.950472	5.567464	0						
50536	19.416488	15.901258	7.950472	7.874008	0						

2194 rows x 5 columns

#Duration										
q1=np.quantile(df_new["Duration"], 0.25)										
q2=np.quantile(df_new["Duration"], 0.50)										
q3=np.quantile(df_new["Duration"], 0.75)										

print(q1,q2,q3)

3.04.69041575982343 7.211102550929798

IQR=q3-q1										
IQR										

4.211102550929798

IQR=q1-1.5*IQR										
IQRq1=1.5*IQR										
print(IQR,IQRq1)										

13.527756377319946 -3.316653826391967

df_new[df_new["Duration"] > IQR]										
Duration Net Sales Age Claim Agency Agency Type Distribution Channel Product Name Destination										
18	15.491933	5.385165	5.099020	0	5	1	0	12	88	

64	19.646883	14.696938	7.280110	0	2	0	1	4	79	
99	19.364917	19.731954	5.099020	0	2	0	1	3	79	

100	13.628182	6.557439	6.000000	0	7	1	1	1	17	
111	20.760539	16.501515	5.830952	1	2	0	1	4	79	

...										
50493	19.157244	12.688578	7.071068	0	2	0	1	4	79	

50513	19.104973	15.901258	7.950472	0	2	0	1	4	79	
50525	19.369506	6.292853	6.403124	0	6	1	1	16	36	

50536	18.416488	15.901258	7.950472	0	2	0	1	4	79	
50546	19.784049	3.162278	6.000000	0	7	1	1	10	88	

2502 rows x 9 columns

df_new[df_new["Duration"] < IQR]										
Duration Net Sales Age Claim Agency Agency Type Distribution Channel Product Name Destination										
1	9.219544	8.306624	6.000000	0	7	1	1	10	79	

2	3.316625	4.449719	8.660254	0	6	1	1	16	56	
3	4.000000	4.472136	5.656854	0	7	1	1	1	38	

4	3.162278	3.872983	5.385165	0	7	1	1	1	47	
5	8.000000	7.035624	6.000000	0	6	1	1	16	88	

...										
50548	3.316625	9.949874	7.071068	0	6	1	1	16	65	

50549	6.928203	4.472136	6.000000	0	7	1	1	1	38	
50550	4.898979	5.291503	6.928203	0	10	1	1	24	38	

50551	5.000000	3.872983	6.000000	0	7	1	1	1	7	
50552	3.000000	4.472136	7.681146	0	7	1	1	0	79	

46034 rows x 9 columns

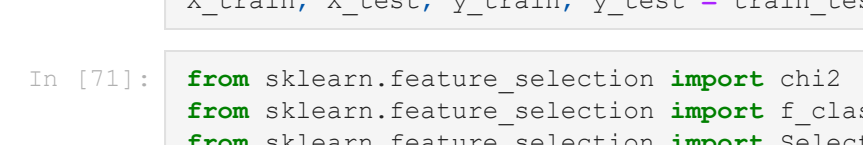
Labels Checking

LABELS = ["Non_claim","Claim"]										
--------------------------------	--	--	--	--	--	--	--	--	--	--

count_classes = pd.value_counts(df_new['Claim'], sort = True)										
count_classes.plot(kind = 'bar', rot=0)										

plt.title("Insurance Claim")										
plt.xticks(range(2), LABELS)										
plt.xlabel("Claim")										
plt.ylabel("Frequency")										

Text(0, 0.5, 'Frequency')



```

fe=SelectKBest(score_func=chi2,k=5)

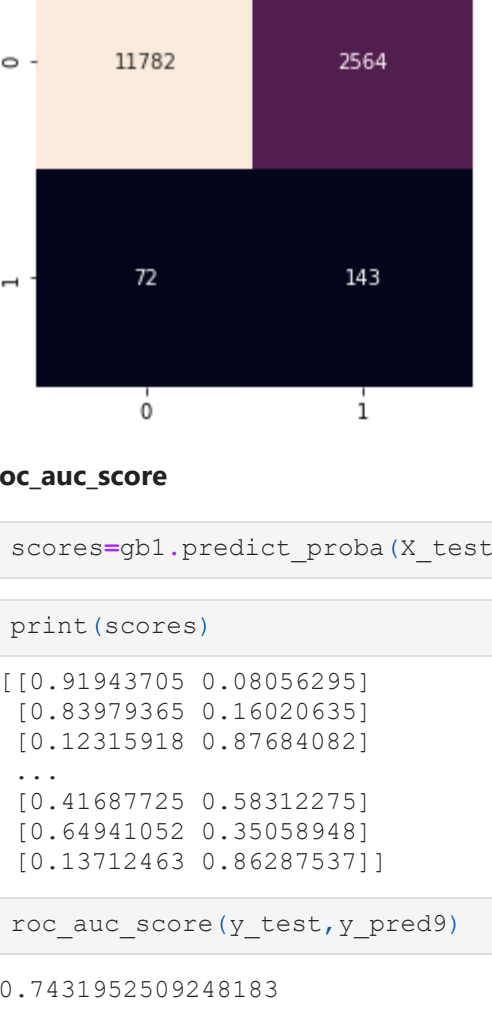
X_train_fe=fe.fit_transform(X_train,y_train)

```


From above fig we cn see our model accuracy with gimli index is 82%

```
In [137]: mat9=confusion_matrix(y_test,y_pred9)
sns.heatmap(mat9,square=True,annot=True,fmt='d')

Out[137]: <AxesSubplot>



roc_auc_score

In [138]: scores=gbl.predict_proba(X_test_ss)

In [139]: print(scores)

[[0.91943705 0.08056295]
 [0.83797865 0.16202135]
 [0.12315918 0.87684082]
 ...
 [0.41687725 0.58312275]
 [0.64941052 0.35058948]
 [0.13712463 0.86287377]]

In [140]: roc_auc_score(y_test,y_pred9)

Out[140]: 0.7431952509248183

Model 5= AdaBoost Classifier

In [141]: from sklearn.ensemble import AdaBoostClassifier
absl=AdaBoostClassifier(n_estimators=100)

In [142]: abal.fit(X_sample2,y_sample2)

Out[142]: AdaBoostClassifier(n_estimators=100)

In [143]: score=cross_val_score(abal,X_sample2,y_sample2,cv=5)
score

Out[143]: array([0.76915029, 0.76606624, 0.76339586, 0.76787983, 0.76674141])

In [144]: y_pred11=abal.predict(X_test_ss)
print(classification_report(y_test,y_pred11))
mat11=confusion_matrix(y_test,y_pred11)
sns.heatmap(mat11,square=True,annot=True,fmt='d')



```

 precision recall f1-score support

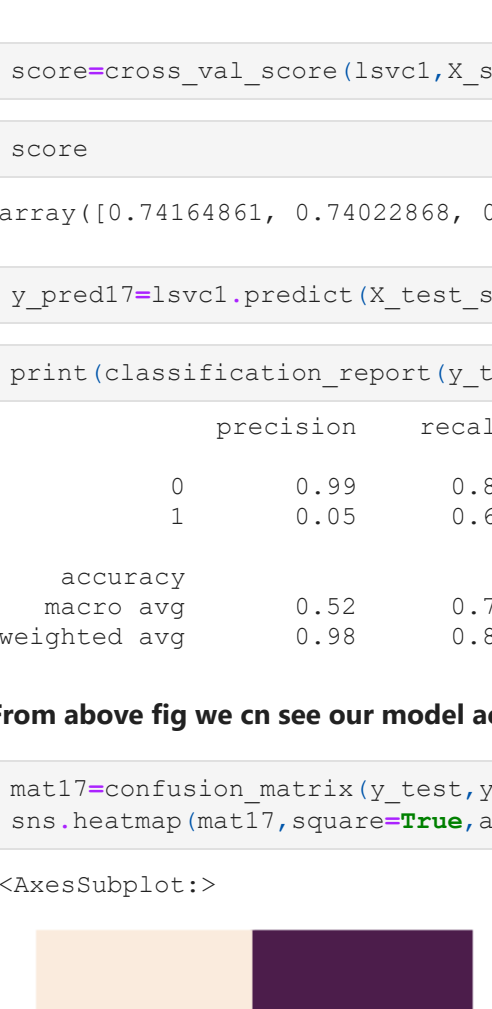
 0 0.99 0.81 0.89 14346
 1 0.05 0.67 0.09 215

 accuracy 0.52 0.74 0.49 14561
 macro avg 0.52 0.74 0.49 14561
 weighted avg 0.98 0.81 0.88 14561

```



Out[144]: <AxesSubplot>


```

From above fig we cn see our model accuracy with gimli index is 82%

```
roc_auc_score

In [145]: scores=abal.predict_proba(X_test_ss)
scores=abal.predict_proba(X_test_ss)

In [146]: print(scores)

[[0.50735607 0.49264393]
 [0.50528633 0.49471367]
 [0.49455375 0.50544625]
 ...
 [0.50002849 0.49997151]
 [0.50190542 0.49809458]
 [0.49492806 0.50507194]]

In [147]: roc_auc_score(y_test,y_pred11)

Out[147]: 0.7374936660928093

Model 6 =SVC

In [148]: lsvc1=LinearSVC()

In [149]: lsvc1.fit(X_sample2,y_sample2)

Out[149]: LinearSVC()

In [150]: score=cross_val_score(lsvc1,X_sample2,y_sample2,cv=5)

In [151]: score

Out[151]: array([0.74164861, 0.74022868, 0.74022868, 0.74336746, 0.74050822])

In [152]: y_pred17=lsvc1.predict(X_test_ss)

In [153]: print(classification_report(y_test,y_pred17))



```

 precision recall f1-score support

 0 0.99 0.83 0.90 14346
 1 0.05 0.67 0.10 215

 accuracy 0.52 0.75 0.50 14561
 macro avg 0.52 0.75 0.50 14561
 weighted avg 0.98 0.83 0.89 14561


```


```

From above fig we cn see our model accuracy with gimli index is 83%

```
In [154]: mat17=confusion_matrix(y_test,y_pred17)
sns.heatmap(mat17,square=True,annot=True,fmt='d')

Out[154]: <AxesSubplot>



Soft Margin

In [155]: lsrv1=LinearSVC(C=0.5,random_state=1)

In [156]: lsrv1.fit(X_sample2,y_sample2)

Out[156]: LinearSVC(C=0.5, random_state=1)

In [157]: score=cross_val_score(lsrv1,X_sample2,y_sample2,cv=5)

In [158]: score

Out[158]: array([0.74164861, 0.74022868, 0.74022868, 0.74336746, 0.74050822])

In [159]: y_pred18=lsrv1.predict(X_test_ss)

In [160]: print(classification_report(y_test,y_pred18))



```

 precision recall f1-score support

 0 0.99 0.83 0.90 14346
 1 0.05 0.67 0.10 215

 accuracy 0.52 0.75 0.50 14561
 macro avg 0.52 0.75 0.50 14561
 weighted avg 0.98 0.83 0.89 14561

```

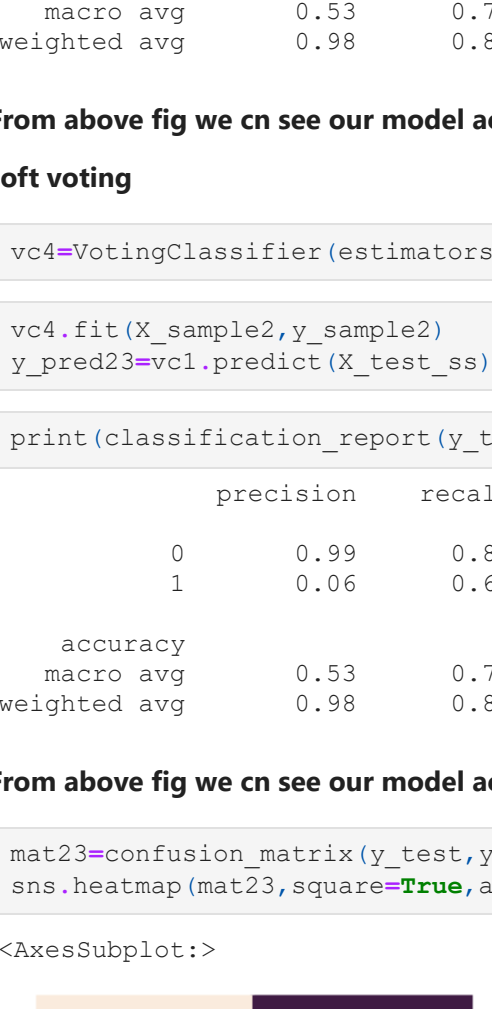

```

From above fig we cn see our model accuracy with gimli index is 82%

```
In [161]: mat18=confusion_matrix(y_test,y_pred18)

In [162]: sns.heatmap(mat18,square=True,annot=True,fmt='d')

Out[162]: <AxesSubplot>


```

Ensembling

Hard voting

```
In [163]: from sklearn.ensemble import VotingClassifier

In [164]: model_list1 = [("lr1",lr1),("dt1",dt1),("dt3",dt3),("rfc1",rfc1),("gbl",gbl),("abal",abal),("lsvc1",lsvc1)]

In [165]: vc1=VotingClassifier(estimators=model_list1)
vc1.fit(X_sample2,y_sample2)

Out[165]: VotingClassifier(estimators=[('lr1', LogisticRegression()),
                                     ('dt1',
                                      DecisionTreeClassifier(max_depth=20,
                                                              min_samples_leaf=35)),
                                     ('dt3',
                                      DecisionTreeClassifier(criterion='entropy',
                                                              max_depth=25,
                                                              min_samples_leaf=30)),
                                     ('rfc1', RandomForestClassifier(criterion='entropy',
                                                                      random_state=1)),
                                     ('gbl', GradientBoostingClassifier()),
                                     ('abal', AdaBoostClassifier(n_estimators=100)),
                                     ('lsvc1', LinearSVC())])

In [166]: vc1=VotingClassifier(estimators=model_list1)

In [167]: vc1.fit(X_sample2,y_sample2)
y_pred22=vc1.predict(X_test_ss)
print(classification_report(y_test,y_pred22))



```

 precision recall f1-score support

 0 0.99 0.85 0.92 14346
 1 0.06 0.62 0.11 215

 accuracy 0.53 0.74 0.51 14561
 macro avg 0.53 0.74 0.51 14561
 weighted avg 0.98 0.85 0.91 14561

```


```

From above fig we cn see our model accuracy with gimli index is 85%

soft voting

```
In [168]: vc4=VotingClassifier(estimators=model_list1,voting="soft")

In [169]: vc4.fit(X_sample2,y_sample2)
y_pred23=vc4.predict(X_test_ss)

In [170]: print(classification_report(y_test,y_pred23))



```

 precision recall f1-score support

 0 0.99 0.85 0.92 14346
 1 0.06 0.62 0.11 215

 accuracy 0.53 0.74 0.51 14561
 macro avg 0.53 0.74 0.51 14561
 weighted avg 0.98 0.85 0.91 14561

```


```

From above fig we cn see our model accuracy with gimli index is 85%

```
In [171]: mat23=confusion_matrix(y_test,y_pred23)
sns.heatmap(mat23,square=True,annot=True,fmt='d')

Out[171]: <AxesSubplot>


```

ConclusionIn this project, I have used different ways of addressing the task with unbalanced data. Like Supervised learning LinearSVC, AdaBoost, GradientBoosting, DecisionTreeClassifier, Logistic Regression and RandomForest. As per their's result I conclude that Random forest is the best algorithm model since it's giving high accuracy than other's. The highest accuracy is 95%

```
In [ ]:
```