

Insurance Claims- Fraud Detection

Problem Statement:

Money has now become the most essential commodity along with food and shelter. It plays important role in every aspect may it be as industry, personal growth.

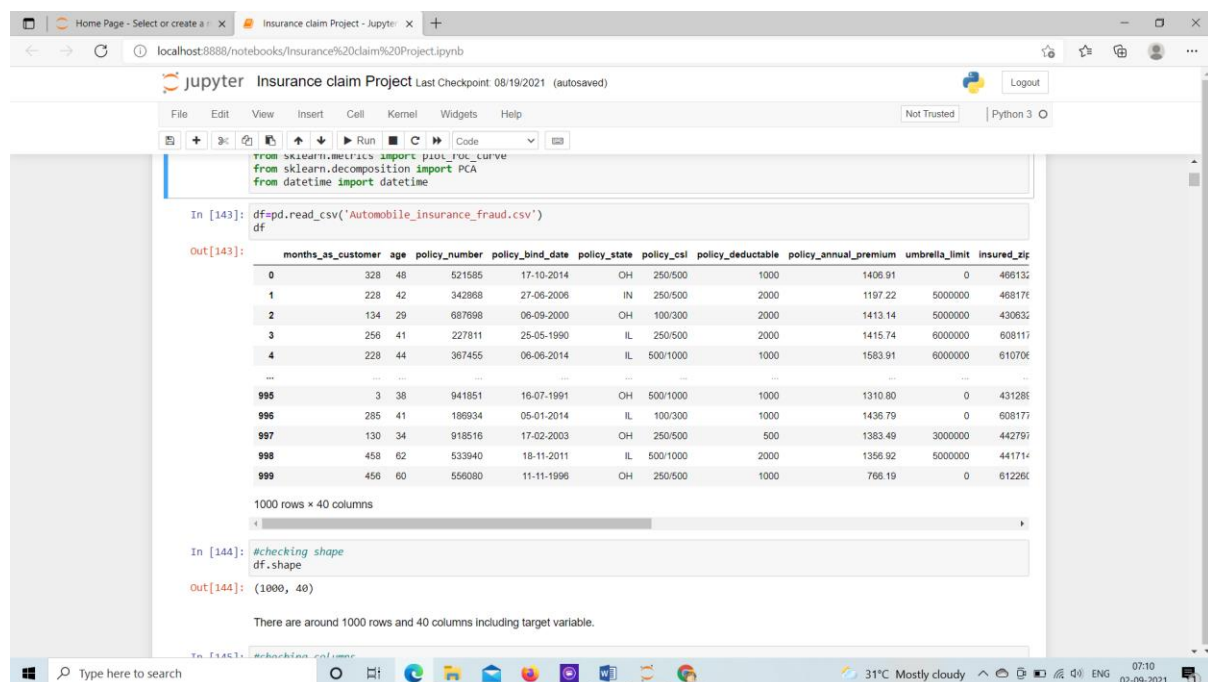
Insurance is also plays an important role in protection from financial loss. It is a form of risk management, primarily used to hedge against the risk of a contingent or uncertain loss.

Insurance fraud is a huge problem in the industry. It's difficult to identify fraud claims. Machine Learning is in a unique position to help the Auto Insurance industry with this problem.

In this project, you are provided a dataset which has the details of the insurance policy along with the customer details. It also has the details of the accident on the basis of which the claims have been made.

In this example, you will be working with some auto insurance data to demonstrate how you can create a predictive model that predicts if an insurance claim is fraudulent or not.

Data Analysis



```
from sklearn.metrics import plot_roc_curve
from sklearn.decomposition import PCA
from datetime import datetime

In [143]: df = pd.read_csv('Automobile_insurance_fraud.csv')
df

Out[143]:
```

	months_as_customer	age	policy_number	policy_bind_date	policy_state	policy_csl	policy_deductable	policy_annual_premium	umbrella_limit	insured_zip
0	328	48	521585	17-10-2014	OH	250/500	1000	1406.91	0	466132
1	228	42	342868	27-06-2006	IN	250/500	2000	1197.22	5000000	468176
2	134	29	687998	06-09-2000	OH	100/300	2000	1413.14	5000000	430632
3	256	41	227811	25-05-1990	IL	250/500	2000	1415.74	6000000	608117
4	228	44	367455	06-06-2014	IL	500/1000	1000	1583.91	6000000	610706
...
995	3	38	941851	16-07-1991	OH	500/1000	1000	1310.80	0	431286
996	285	41	186934	05-01-2014	IL	100/300	1000	1436.79	0	608177
997	130	34	918516	17-02-2003	OH	250/500	500	1383.49	3000000	442797
998	458	62	533940	18-11-2011	IL	500/1000	2000	1356.92	5000000	441714
999	456	60	556080	11-11-1996	OH	250/500	1000	766.19	0	612262

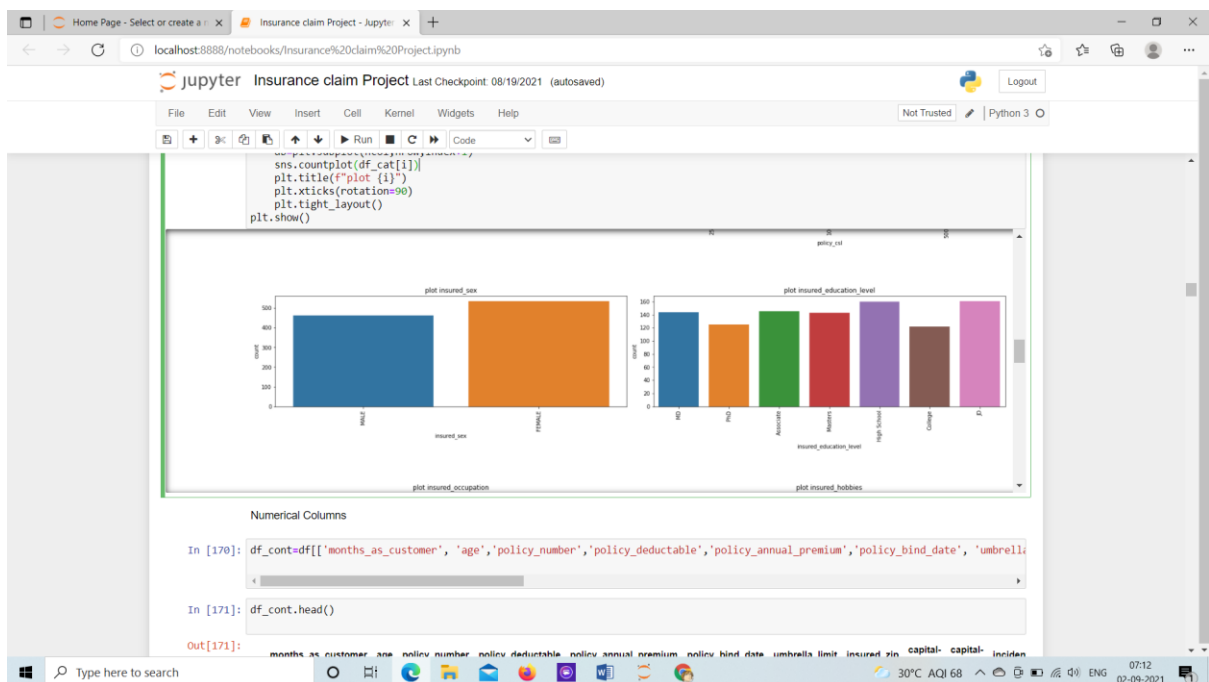
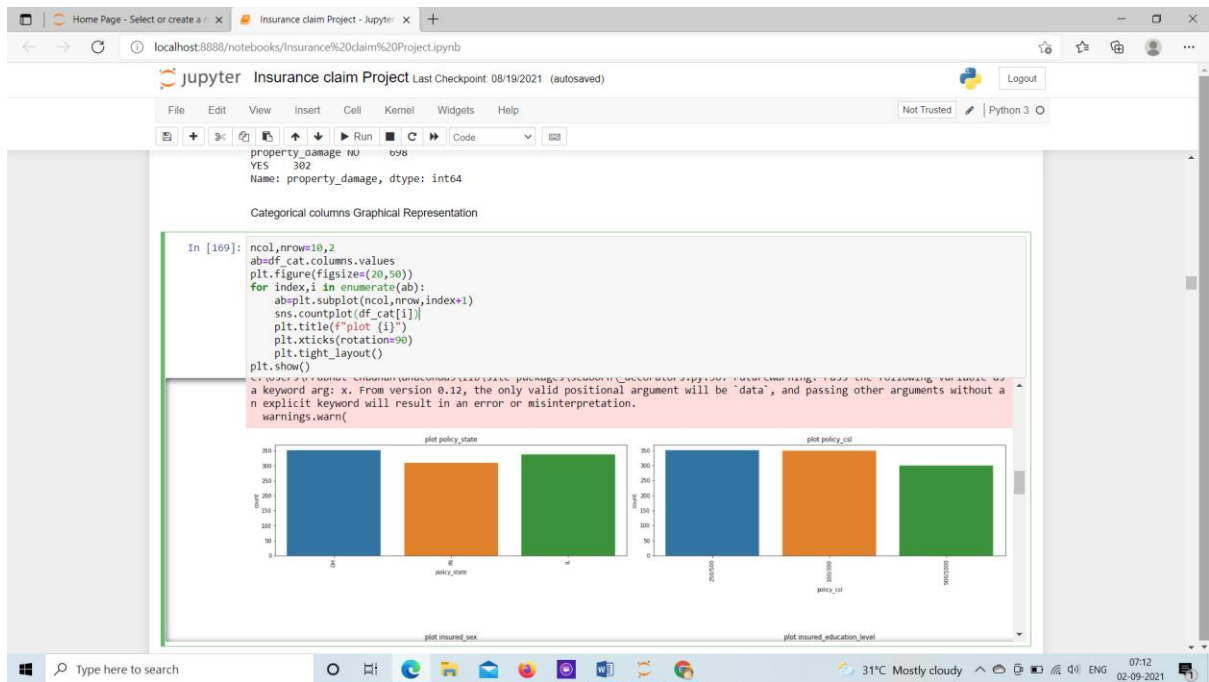
```
1000 rows x 40 columns

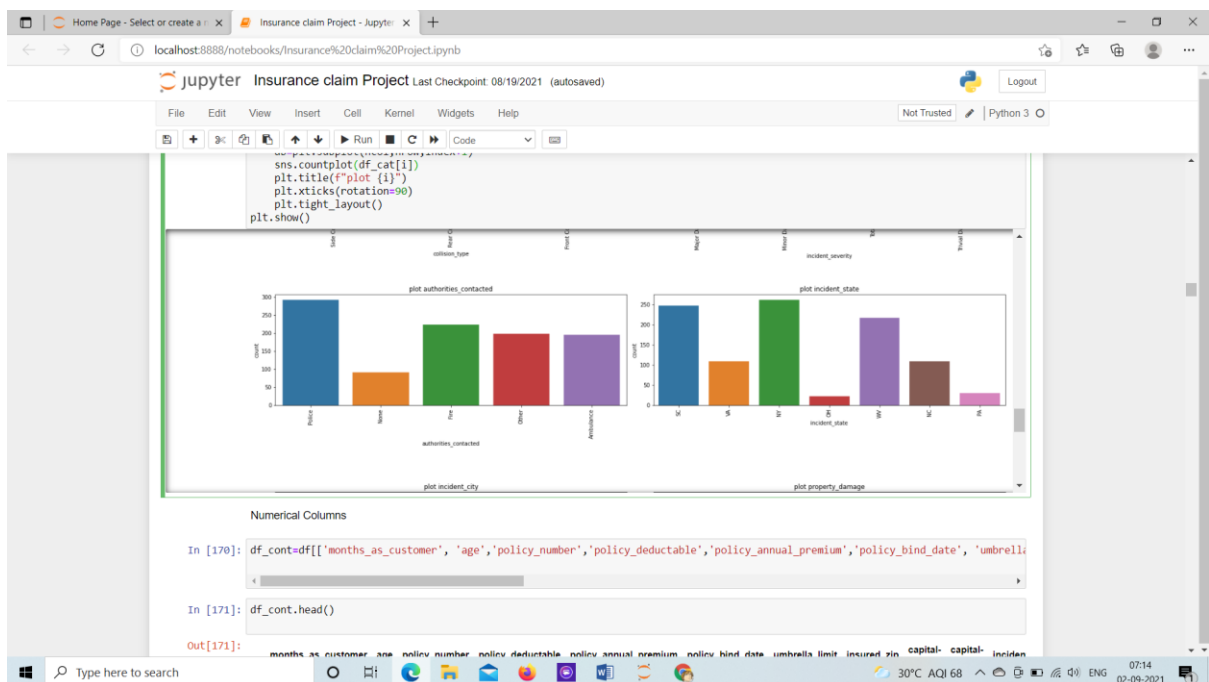
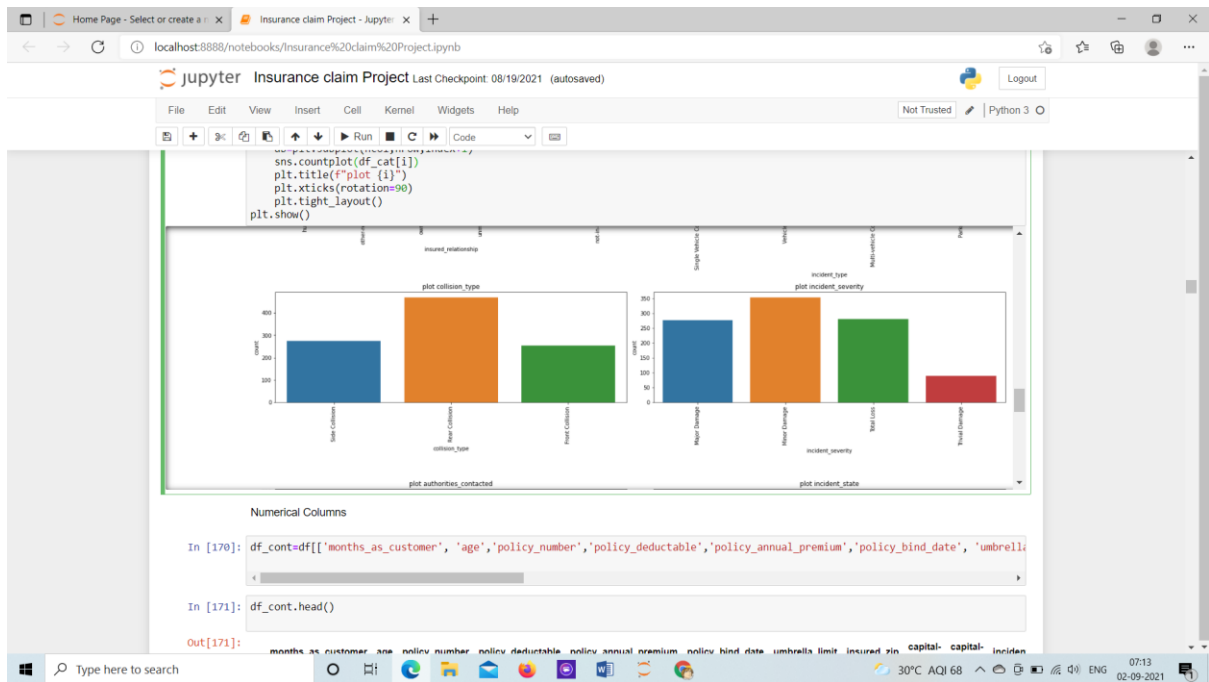
In [144]: #checking shape
df.shape

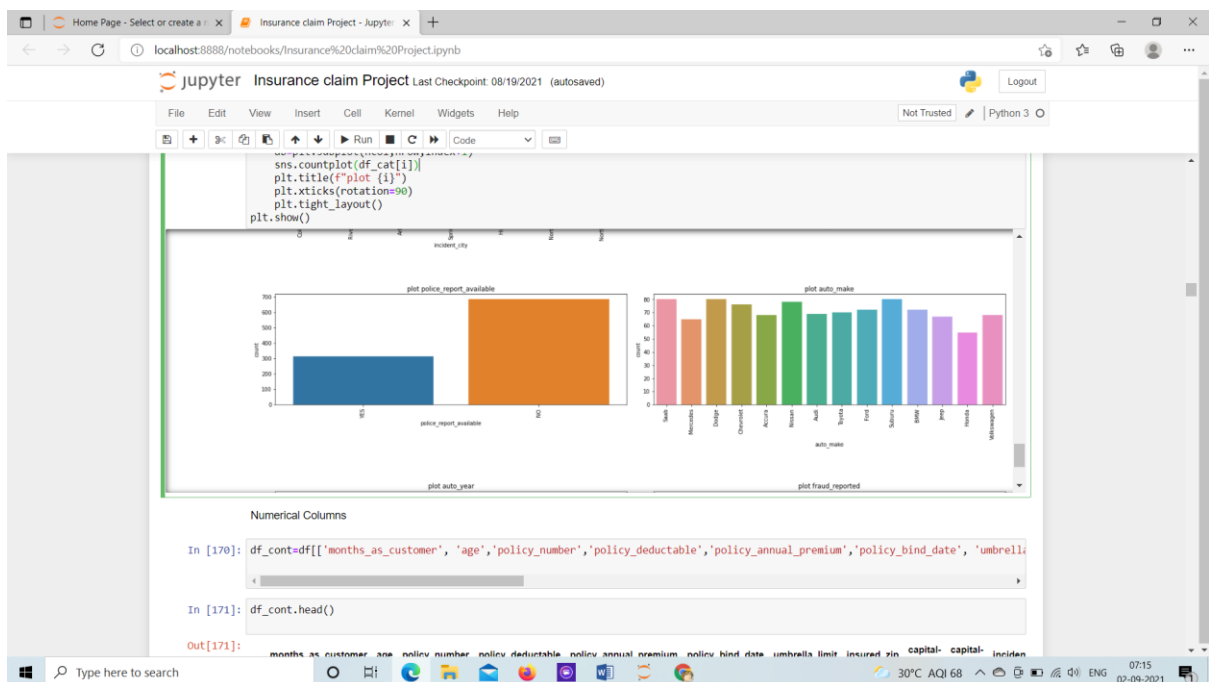
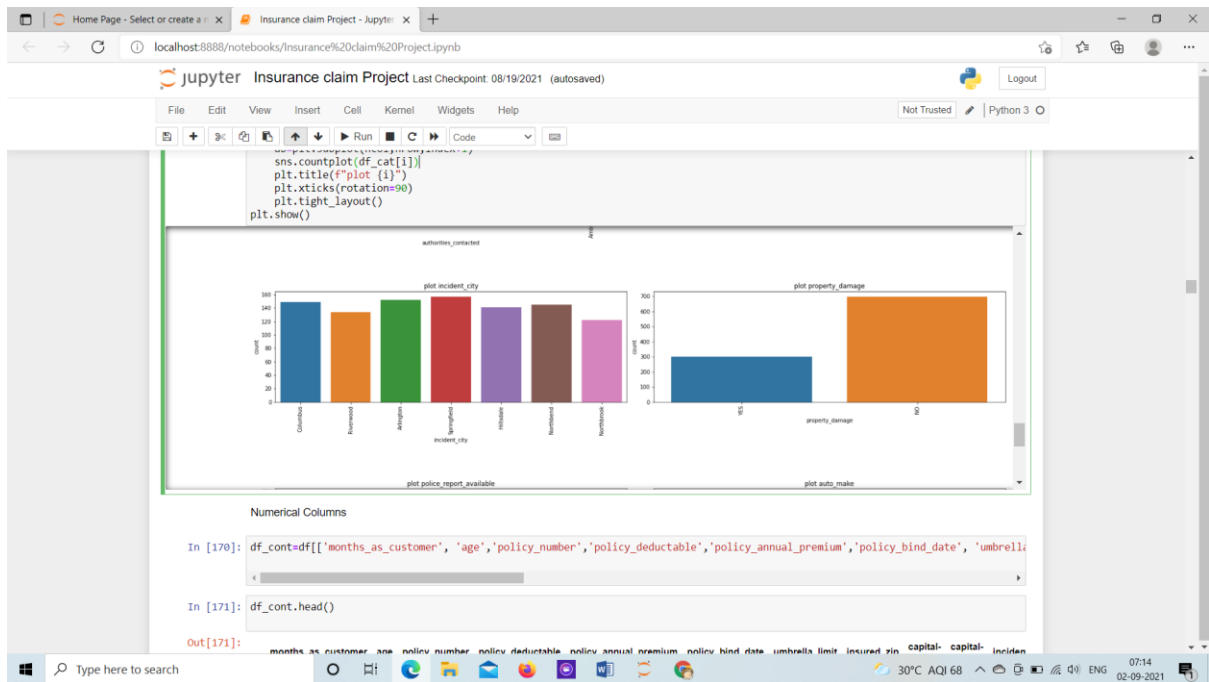
Out[144]: (1000, 40)

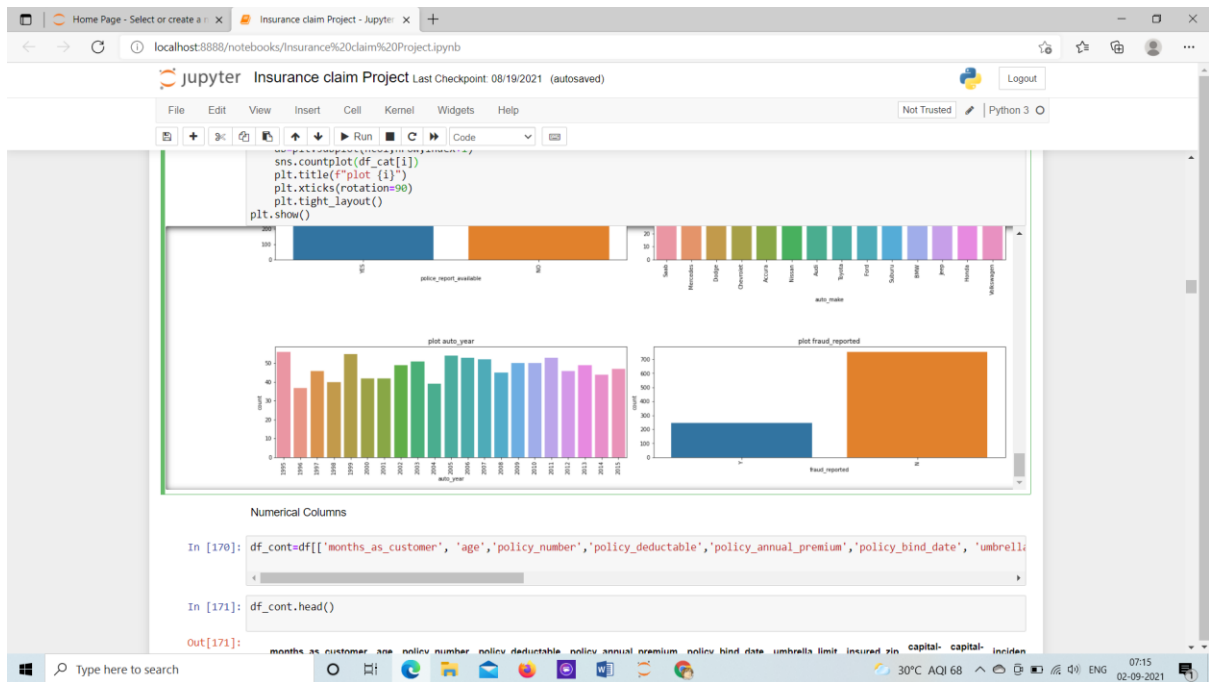
There are around 1000 rows and 40 columns including target variable.
```

EDA

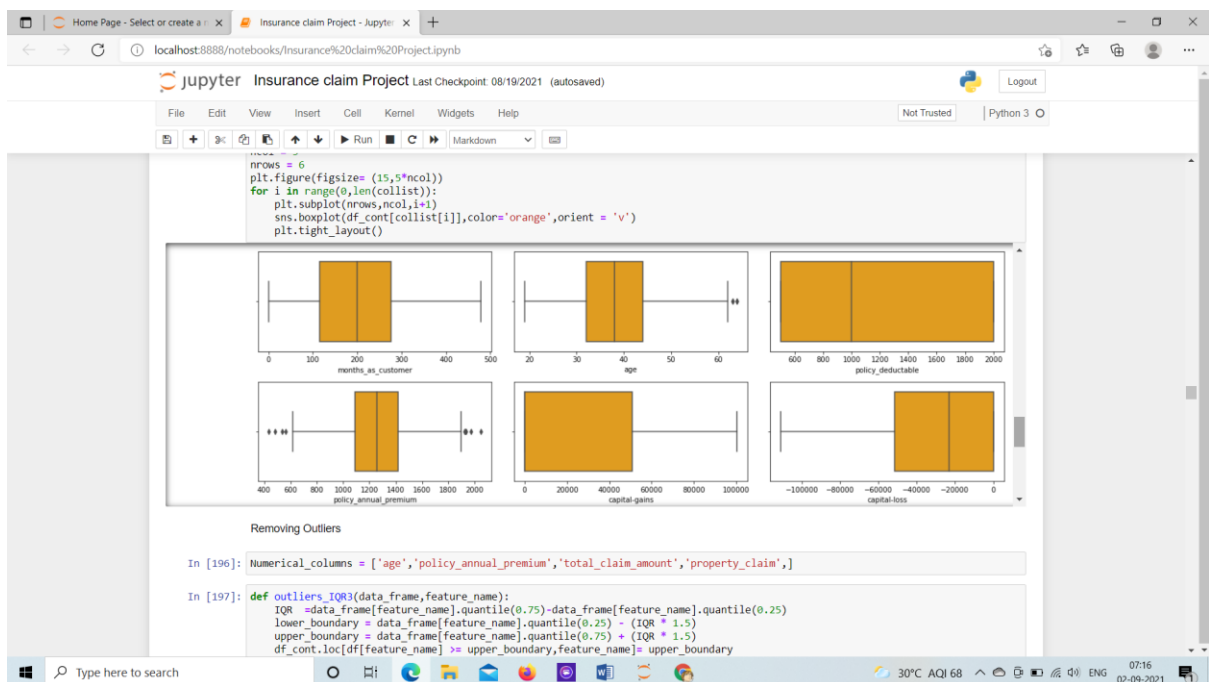


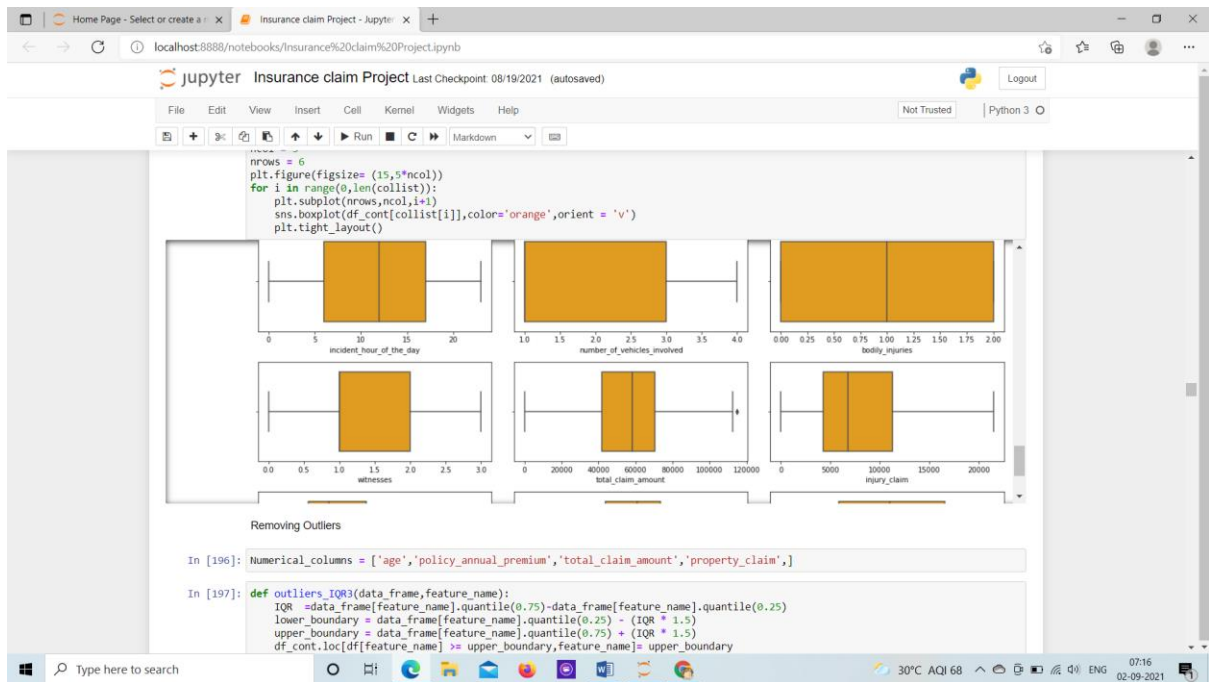




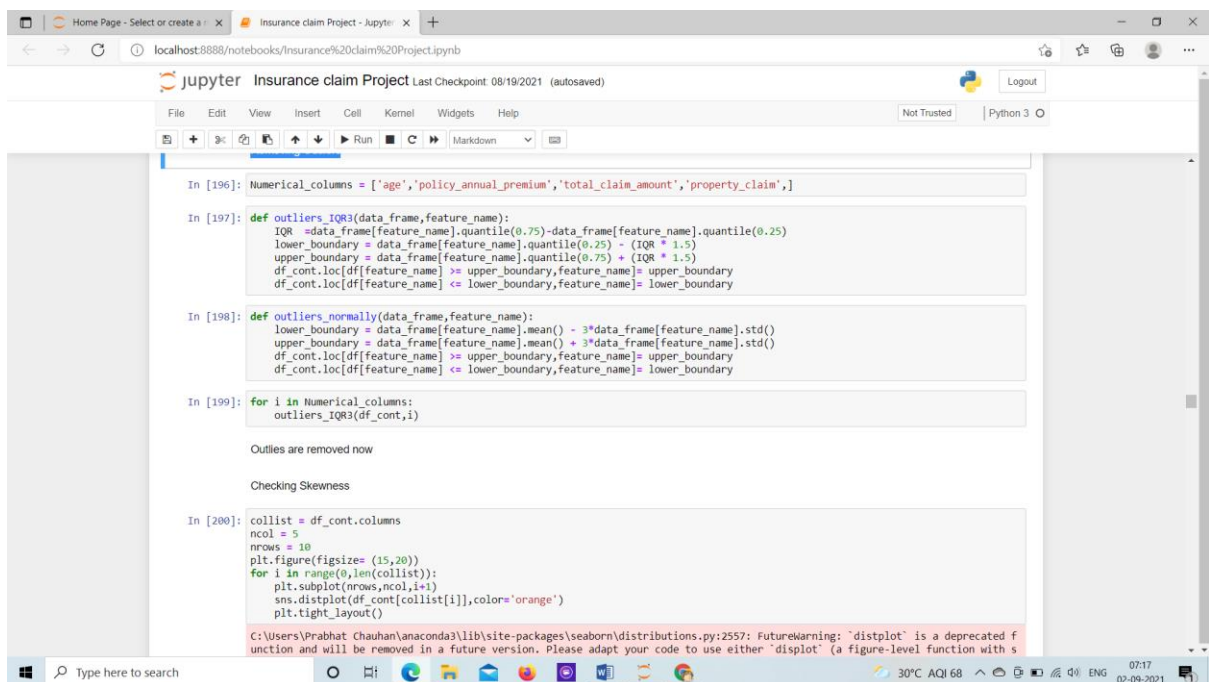


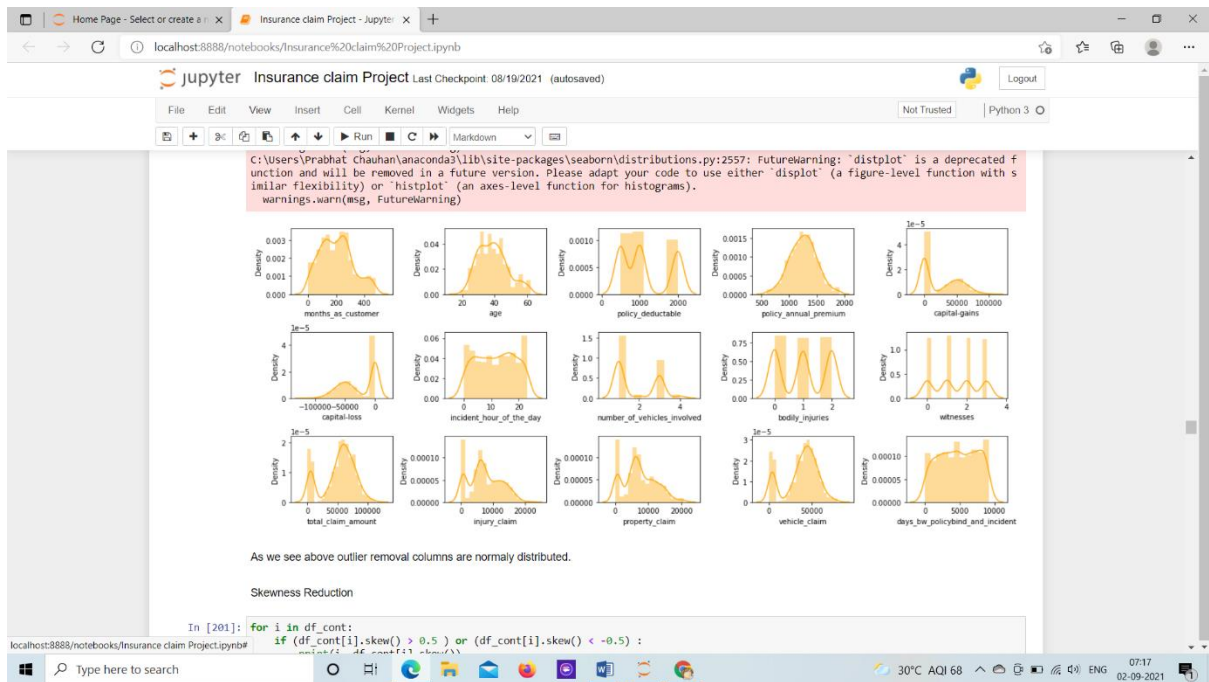
Checking Outliers using Boxplot





Removing Outliers





Random State

```

x_sm, y_sm = smote.fit_resample(pca_x, y)
y_sm.value_counts()

Out[226]: 0.0    753
          1.0    753
          Name: fraud_reported, dtype: int64

Random State

In [227]: maxAccu=0
          maxRS=0
          for i in range(1,200):
              x_train,x_test,y_train,y_test=train_test_split(x_sm,y_sm,test_size=.30,random_state=i,stratify= y_sm)
              LR=LogisticRegression()
              LR.fit(x_train,y_train)
              pred=LR.predict(x_test)
              acc=accuracy_score(y_test,pred)
              if acc>maxAccu:
                  maxAccu=acc
                  maxRS=i
          print("Best accuracy is",maxAccu,"n Random_state",maxRS)

Best accuracy is 0.8805309734513275 n Random_state 10

Model

In [228]: x_train,x_test,y_train,y_test = train_test_split(x_sm,y_sm,test_size=.33,random_state= 115, stratify= y_sm )

In [229]: y_train.value_counts()

Out[229]: 1.0    505
          0.0    504
          Name: fraud_reported, dtype: int64

```

Model selection


```
0.0      504
Name: fraud_reported, dtype: int64

Model selection

In [230]: models=[LogisticRegression(),
                  SVC(),
                  GaussianNB(),
                  DecisionTreeClassifier(),
                  KNeighborsClassifier(),
                  RandomForestClassifier(),
                  AdaBoostClassifier(),
                  GradientBoostingClassifier(),
                  BaggingClassifier()]
least_difference = []
for m in models:
    print("\n")
    print(m)
    m.fit(x_train,y_train)
    pred = m.predict(x_test)
    accu = accuracy_score(y_test,pred)
    f1 = f1_score(y_test,pred)
    print(confusion_matrix(y_test,pred))
    print(classification_report(y_test,pred))
    print("The accuracy of {} is {}".format(m,accu))
    cv = cross_val_score(m,x_test,y_test,cv=5)
    print("cross val score :", cv.mean())
    difference = np.abs(accuracy_score(y_test,pred) - cv.mean())
    a = 'difference b/w accuracy score and cross val score is : {:.2f}'.format(difference)
    print(a)
    least_difference.append((m,a))
    for i in [0,1]:
        f1 = f1_score(y_test,pred,pos_label=i)
        print("f1 score for {} is {}".format(i,f1))

f1 score for 0 is 0.8843813387423934
f1 score for 1 is 0.8862275449101796
```

```
print("f1 score for {} is {}".format(1,f1))

LogisticRegression()
[[220  29]
 [ 33 215]]
      precision    recall  f1-score   support

 0.0       0.87       0.88       0.88        249
 1.0       0.88       0.87       0.87        248

 accuracy          0.88          497
 macro avg         0.88          497
 weighted avg      0.88          497

The accuracy of LogisticRegression() is 0.8752515090543259
cross val score : 0.8539394072737675
difference b/w accuracy score and cross val score is : 0.02
f1 score for 0 is 0.8764940239043825
f1 score for 1 is 0.8739837398373985

Finding the least difference

In [231]: least_difference
Out[231]: [(LogisticRegression(),
'difference b/w accuracy score and cross val score is : 0.02'),
(SVC(), 'difference b/w accuracy score and cross val score is : 0.01'),
(GaussianNB(),
'difference b/w accuracy score and cross val score is : 0.01'),
(DecisionTreeClassifier(),
'difference b/w accuracy score and cross val score is : 0.00'),
(KNeighborsClassifier(),
'difference b/w accuracy score and cross val score is : 0.01'),
(RandomForestClassifier(),
'difference b/w accuracy score and cross val score is : 0.02'),
(AdaBoostClassifier(),
'difference b/w accuracy score and cross val score is : 0.02')]
```

Hyper Tuning Model

```

In [232]: KNN_parameters={'n_neighbors':np.arange(1, 16),
'weights':{'uniform', 'distance'},
'algorithm':{'auto','ball_tree','kd_tree','brute'}}

In [233]: KNN=GridSearchCV(KNeighborsClassifier(),KNN_parameters,cv=5)

In [234]: KNN.fit(x_train,y_train)

Out[234]: GridSearchCV(cv=5, estimator=KNeighborsClassifier(),
param_grid={'algorithm': ('auto', 'ball_tree', 'kd_tree', 'brute'),
'n_neighbors': array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15]),
'weights': ('uniform', 'distance')})

In [235]: KNN.best_params_

Out[235]: {'algorithm': 'auto', 'n_neighbors': 2, 'weights': 'uniform'}

In [236]: final_model =KNeighborsClassifier(algorithm='auto', n_neighbors= 2,weights='uniform')
final_model.fit(x_train,y_train)
pred = final_model.predict(x_test)
accu =accuracy_score(y_test,pred)
print(confusion_matrix(y_test,pred))
print(classification_report(y_test,pred))
print('The accuracy of {} is {}'.format(final_model,accu))
cv = cross_val_score(final_model,x_sm,y_sm,cv =5)
print('\n')
print("cross val score :", cv.mean())
difference = np.abs(accuracy_score(y_test,pred) - cv.mean())
least_difference.append((final_model,difference))
print('\n')
print('difference b/w accuracy score and cross val score is :',difference)
for i in ([0,1]):
    f1 = f1_score(y_test,pred,pos_label=i)
    print("f1 score for {} is {}".format(i,f1))

```

```

print('\n')
print("cross val score :", cv.mean())
difference = np.abs(accuracy_score(y_test,pred) - cv.mean())
least_difference.append((final_model,difference))
print('\n')
print('difference b/w accuracy score and cross val score is :',difference)
for i in ([0,1]):
    f1 = f1_score(y_test,pred,pos_label=i)
    print("f1 score for {} is {}".format(i,f1))

[[139 110]
 [ 10 238]]
precision    recall  f1-score   support

   0.0         0.93    0.56    0.70     249
   1.0         0.68    0.96    0.80     248

 accuracy
macro avg    0.81    0.76    0.75     497
weighted avg    0.81    0.76    0.75     497

The accuracy of KNeighborsClassifier(n_neighbors=2) is 0.7585513078470825

cross val score : 0.7676068733361202

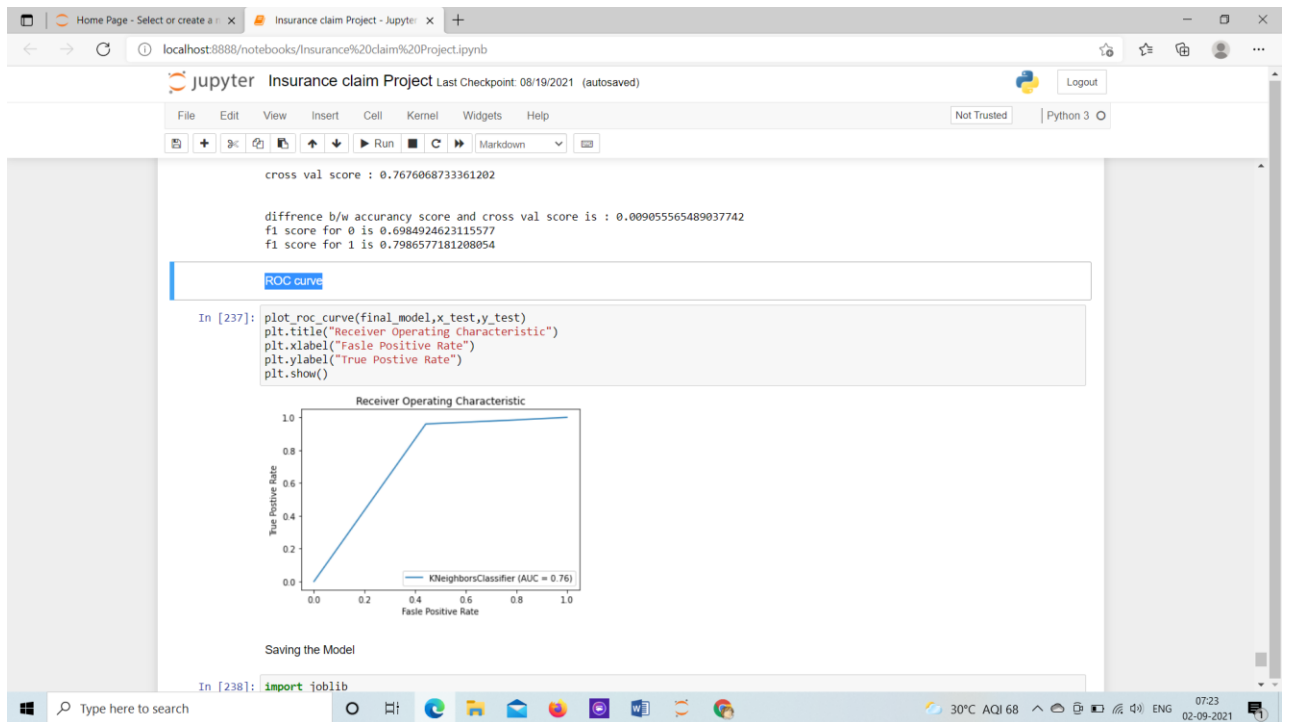
difference b/w accuracy score and cross val score is : 0.009055565489037742
f1 score for 0 is 0.6984924623115577
f1 score for 1 is 0.7986577181208054

ROC curve

In [237]: plot_roc_curve(final_model,x_test,y_test)
plt.title("Receiver Operating characteristic")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")

```

ROC curve



Result

Accuracy proportion we have is 75%

