# Optimizing Flight Booking Decisions through Machine Learning Price Predictions
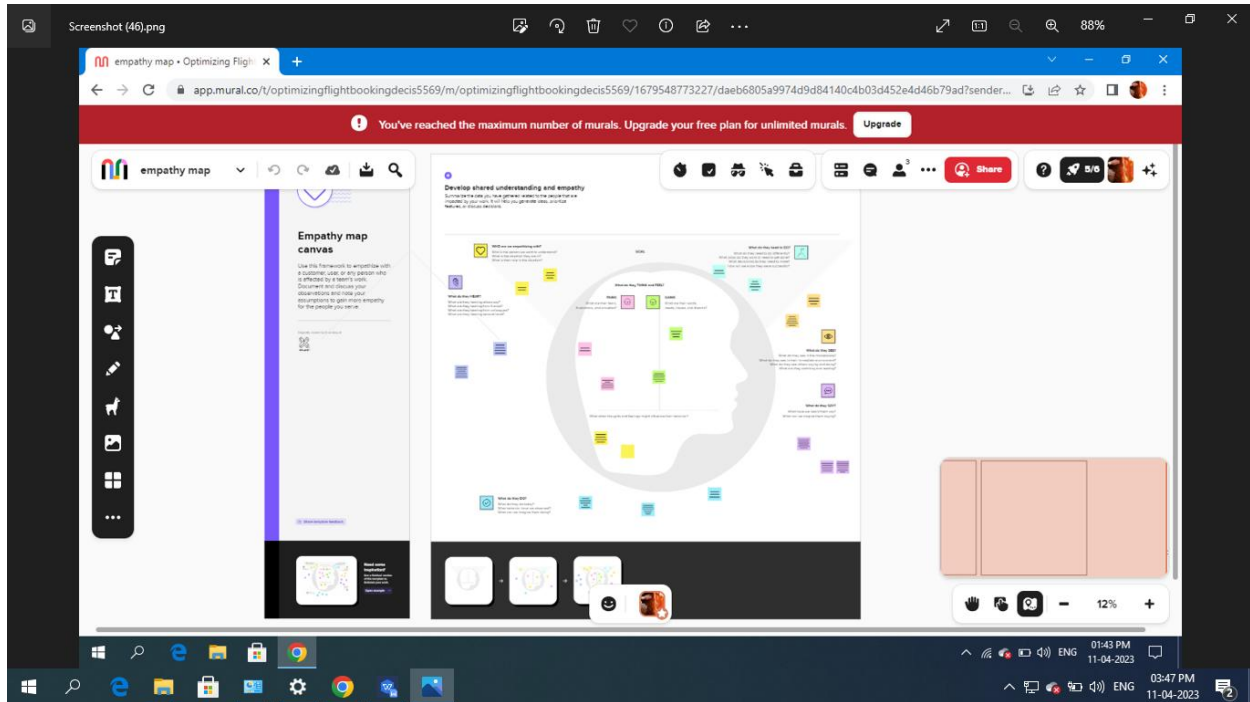
## OVERVIEW:

In this article, we will be **analyzing the flight fare prediction using Machine Learning dataset** using essential exploratory data analysis techniques then will **draw some predictions about the price of the flight based on some features** such as what type of airline it is, what is the arrival time, what is the departure time, what is the duration of the flight, source, destination and more.
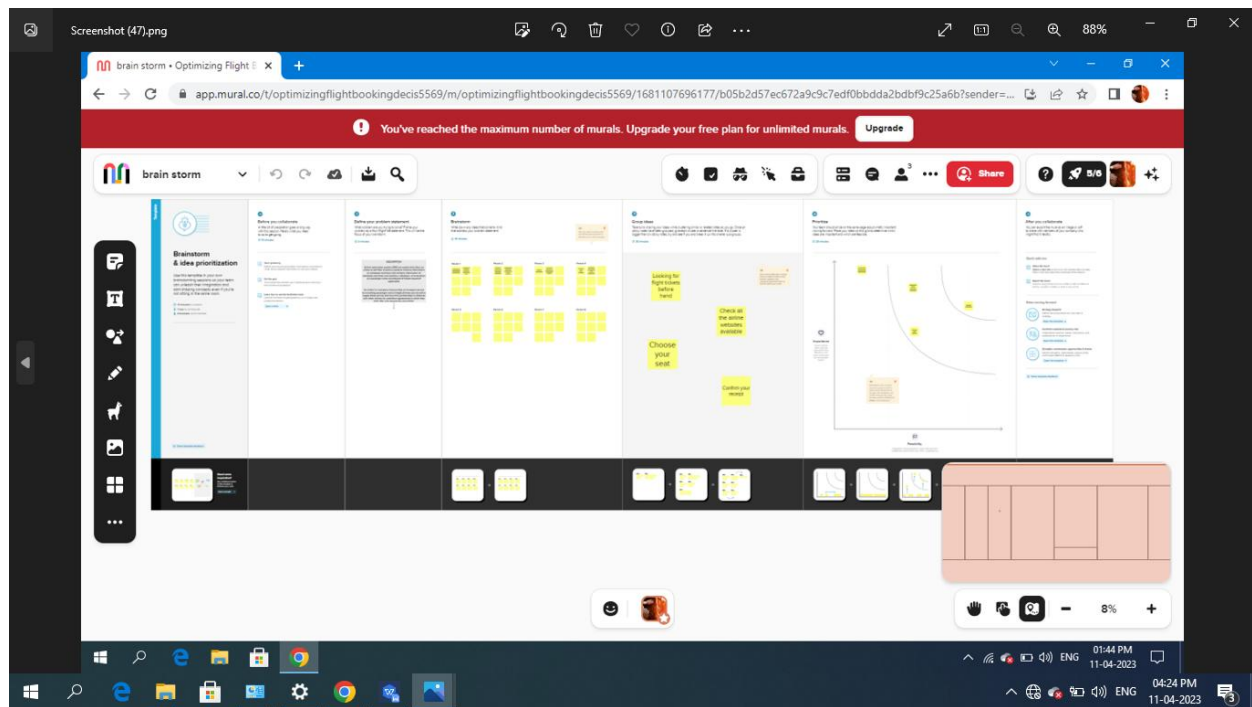
## PURPOSE:

 **Optimizing  Flight Booking Decisions through Machine Learning Price Predictions to provide convenient to passenger.**

# Problem Definition & Design Thinking

## Empathy Map:

# Ideation & Brainstorming Map:

# RESULT:

**Screenshot 1 — flight booking.ipynb**

```python
data.isnull().sum()
```

```
Airline               0
Date_of_Journey       0
Source                0
Destination           0
Route                 1
Dep_Time              0
Arrival_Time          0
Duration              0
Total_Stops           1
Additional_Info       0
Price                 0
Date                  0
Month                 0
Year                  0
City1                 1
City2             10683
City3             10683
City4             10683
City5             10683
City6             10683
Dep_Time_Hour         0
Dep_Time_Mins         0
Arrival_date       6348
Time_of_Arrival       0
Arrival_Time_Hour     0
Arrival_Time_Mins     0
Travel_Hours          0
```

completed at 12:24 PM

Flight.py

```python
47 data['Time_of_Arrival']=data.Arrival_Time.str
48 data['Time_of_Arrival']=data.Time_of_Arrival.
49 data['Arrival_Time_Hour']=data.Time_of_Arriva
50 data['Arrival_Time_Mins']=data.Time_of_Arriva
51
52 data.Duration=data.Duration.str.split(' ')
53
54 data['Travel_Hours']=data.Duration.str[0]
55 data['Travel_Hours']=data['Travel_Hours'].str
56 data['Travel_Hours']=data['Travel_Hours'].str
57 data.Travel_Hours=data.Travel_Hours
58 data['Travel_Mins']=data.Duration.str[1]
59 data.Travel_Mins=data.Travel_Mins.str.split('
60 data.Travel_Mins=data.Duration.str[0]
61
62 data.Total_Stops.replace('non_stop',0,inplace
63 data.Total_Stops=data.Total_Stops.str.split('
64 data.Total_Stops=data.Total_Stops.str[0]
65
66 data.Additional_Info.unique()
67
68 data.Additional_Info.replace('No Info','No In
69
70 data.isnull().sum()
71
72 data.drop(['City4','City5','City6'],axis=1,in
```

---

**Screenshot 2 — flight booking.ipynb**

```python
data.isnull().sum()
```

```
Airline               0
Source                0
Destination           0
Total_Stops           1
Additional_Info       0
Price                 0
Date                  0
Month                 0
Year                  0
City1                 1
City2             10683
City3             10683
Dep_Time_Hour         0
Dep_Time_Mins         0
Arrival_date       6348
Arrival_Time_Hour     0
Arrival_Time_Mins     0
Travel_Hours          0
Travel_Mins           0
dtype: int64
```

completed at 12:27 PM

Flight.py

```python
61
62 data.Total_Stops.replace('non_stop',0,inplace
63 data.Total_Stops=data.Total_Stops.str.split('
64 data.Total_Stops=data.Total_Stops.str[0]
65
66 data.Additional_Info.unique()
67
68 data.Additional_Info.replace('No Info','No In
69
70 data.isnull().sum()
71
72 data.drop(['City4','City5','City6'],axis=1,in
73 data.drop(['Date_of_Journey','Route','Dep_Tim
74 data.drop(['Time_of_Arrival'],axis=1,inplace=
75
76 data.isnull().sum()
77
78 data['City3'].fillna('None,inplae=True')
79 data['Arrival_Date'].fillna(data['Datedata.De
80 data.Dep_Time_Mins=data.Dep_Time_Mins.astype(
81 data.Arrival_date=data.Arrival_date.astype('i
82 data.Arrival_Time_Hours=data.Arrival_Time_Hou
83 data.Arrival_Time_Mins=data.Arrival_Time_Mins
84 data.Travel_Mins=data.Travel_Mins.astype('int
85
86 data[data['Travel_Hours']=='5m']],inplace=Tr
```

flight booking.ipynb
File Edit View Insert Runtime Tools Help  All changes saved

Comment  Share

Files

+ Code  + Text

```
#Distribution of 'PRICE' Column
plt.figure(figsize=(15,8))
sns.distplot(data.Price)
```

```
<Axes: xlabel='Price', ylabel='Density'>
```



Flight.py

```
115 data.head()
116
117
118
119 import seaborn as sns
120 c=1
121 plt.figure(figsize=(20,45))
122
123 for i in categorical:
124     plt.subplot(6,3,c)
125     sns.countplot(data[i])
126     plt.xticks(rotation=90)
127     plt.tight_layout(pad=3.0)
128     c=c+1
129
130 plt.show()
131
132
133
134 #Distribution of 'PRICE' Column
135 plt.figure(figsize=(15,8))
136 sns.distplot(data.Price)
137
138 sns.heatmap(data.corr(annot=True))
139
140 #Detecting the Outliers
```

1s  completed at 12:41 PM

Disk  84.49 GB available

---

flight booking.ipynb
File Edit View Insert Runtime Tools Help  All changes saved

Comment  Share

Files

+ Code  + Text

```
#Detecting the Outliers
import seaborn as sns
sns.boxplot(data['Price'])
```

```
<Axes: >
```



Flight.py

```
129
130 plt.show()
131
132
133
134 #Distribution of 'PRICE' Column
135 plt.figure(figsize=(15,8))
136 sns.distplot(data.Price)
137
138 sns.heatmap(data.corr(annot=True))
139
140 #Detecting the Outliers
141 import seaborn as sns
142 sns.boxplot(data['Price'])
143
144 y=data['Price']
145 x=data.drop(columns=['Price'],axis=1)
146
147 from sklearn.preprocessing import StandardSca
148 ss=StandardScaler()
149
150 x_scaled=ss.fit_transform
151
152 x_scaled=pd.DataFrame(x_scaled,columns=x.colu
153 x_scaled.head()
154
```

0s  completed at 12:43 PM

Disk  84.49 GB available

**flight booking.ipynb**
File  Edit  View  Insert  Runtime  Tools  Help  All changes saved

+ Code  + Text

SEARCH STACK OVERFLOW

```python
[32] categorical=['Airline','Source','Destination','Additional_Info','City1']
     numerical=['Total_Stops','Date','Month','Year','Dep_Time_Hour','Dep_Time_Min','A
```

Double-click (or enter) to edit

Flight.py ×

```
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101  _Time_Mins','Travel_Hours','Travel_Mins']
102
103
104
105
106
107
108
```

0s  completed at 12:35 PM

---

**Untitled1.ipynb**
File  Edit  View  Insert  Runtime  Tools  Help

+ Code  + Text

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import f1_score
from sklearn.metrics import classification_report,confusion_matrix
import warnings
import pickle
from scipy import stats
import warnings
warnings.filterwarnings('ignore')
plt.style.use('fivethirtyeight')
data=pd.read_csv("Data_Train.csv")
data.head()
```

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time |
|---|---------|-----------------|--------|-------------|-------|----------|--------------|
| 0 | IndiGo | 24/03/2019 | Banglore | New Delhi | BLR ? DEL | 22:20 | 01:10 22 Mar |
| | | | | | CCU | | |

Flight.py ×

```python
1
2  import numpy as np
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  from sklearn.model_selection import train_tes
7  #from sklearn.ensemble import RandomForestCla
8  from sklearn.tree import DecisionTreeClassifi
9  from sklearn.neighbors import KNeighborsClass
10 from sklearn.metrics import f1_score
11 from sklearn.metrics import classification_re
12 import warnings
13 import pickle
14 from scipy import stats
15 warnings.filterwarnings('ignore')
16 plt.style.use('fivethirtyeight')
17
18 data=pd.read_csv(r"C:\Users\DEEPANAYAKI\Deskt
19 data.head()
20
21 data.Date_of_Journey=data.Date_of_Journey.str
22 data.Date_of_Journey
23
24 for i in category:
25   print(i,data[i].unique())
26 data['Date']=data.Date_of_Journey.str[0]
```

0s  completed at 12:05 PM

# ADVANTAGES :

Booking airline flights early can give passengers cheaper deals for travel because it gives them time to find and track the best price.

It can also allow them to begin planning their trip early.

While this may not appear to be a huge benefit for airports, cheaper airline tickets can lead to increased flight travel.

# DISADVANTAGES:

You need internet access. Reliable internet access is required to check reservations and add bookings that are made over the phone.

You need to be ready for an influx of new customers.

Not all online booking systems are created equal.

# APPLICATION:

Flight booking applications help the airline industry automate the booking process. Users worldwide can book flights on the go using the simple apps, which include features such as quick flight search, download tickets, check and modify booking details, one-tap check-in, and many more

# CONCLUSION:

Data visualization as well so after these steps one can go for the prediction using machine learning model making steps.

# FUTURE SCOPE:

Airline reservation system make the life of passengers very easy as they don't need to stand in queues for getting their seats reserved and they can easily make reservations on any airline just from a single system.

# APPENDIX:

```python
import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

from sklearn.model_selection import train_test_split

#from sklearn.ensemble import RandomForestclassifier,GradientBoosting

from sklearn.tree import DecisionTreeClassifier

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import f1_score

from sklearn.metrics import classification_report,confusion_matrix

import warnings

import pickle

from scipy import stats

warnings.filterwarnings('ignore')

plt.style.use('fivethirtyeight')


data=pd.read_csv(r"C:\Users\DEEPANAYAKI\Desktop\New folder\Data_Train.csv")

data.head()


data.Date_of_Journey=data.Date_of_Journey.str.split('/')

data.Date_of_Journey


for i in category:

  print(i,data[i].unique())

data['Date']=data.Date_of_Journey.str[0]
```

```python
data['Month']=data.Date_of_Journey.str[1]

data['Year']=data.Date_of_Journey.str[2]

data.Total_Stops.unique()

data.Route=data.Route.str.split('->')

data.Route


data['City1']=data.Route.str[0]

data['City2']=data.Route.str[1]

data['City3']=data.Route.str[2]

data['City4']=data.Route.str[3]

data['City5']=data.Route.str[4]

data['City6']=data.Route.str[5]


data.Dep_Time=data.Dep_Time.str.split(':')

data['Dep_Time_Hour']=data.Dep_Time.str[0]

data['Dep_Time_Mins']=data.Dep_Time.str[1]


data.Arrival_Time=data.Arrival_Time.str.split(' ')


data['Arrival_date']=data.Arrival_Time.str[1]

data['Time_of_Arrival']=data.Arrival_Time.str[0]

data['Time_of_Arrival']=data.Time_of_Arrival.str.split(':')

data['Arrival_Time_Hour']=data.Time_of_Arrival.str[0]

data['Arrival_Time_Mins']=data.Time_of_Arrival.str[1]


data.Duration=data.Duration.str.split(' ')
```

```python
data['Travel_Hours']=data.Duration.str[0]

data['Travel_Hours']=data['Travel_Hours'].str.split('h')

data['Travel_Hours']=data['Travel_Hours'].str[0]

data.Travel_Hours=data.Travel_Hours

data['Travel_Mins']=data.Duration.str[1]

data.Travel_Mins=data.Travel_Mins.str.split('m')

data.Travel_Mins=data.Duration.str[0]


data.Total_Stops.replace('non_stop',0,inplace=True)

data.Total_Stops=data.Total_Stops.str.split(' ')

data.Total_Stops=data.Total_Stops.str[0]


data.Additional_Info.unique()


data.Additional_Info.replace('No Info','No Info',inplace=True)


data.isnull().sum()


data.drop(['City4','City5','City6'],axis=1,inplace=True)

data.drop(['Date_of_Journey','Route','Dep_Time','Arrival_Time','Duration'],axis=1,
inplace=True)

data.drop(['Time_of_Arrival'],axis=1,inplace=True)


data.isnull().sum()


data['City3'].fillna('None,inplae=True')
```

```python
data['Arrival_Date'].fillna(data['Datedata.Dep_Time_Hur=data.Dep_Time_Hour.astype
('int64')

data.Dep_Time_Mins=data.Dep_Time_Mins.astype('int64')

data.Arrival_date=data.Arrival_date.astype('int64')

data.Arrival_Time_Hours=data.Arrival_Time_Hours.astype('int64')

data.Arrival_Time_Mins=data.Arrival_Time_Mins.astype('int64')

data.Travel_Mins=data.Travel_Mins.astype('int64')


data[data['Travel_Hours']=='5m']'],inplace=True)

data['Travel_Mins'].fillna(0,inplace=True)


data.info()


#data.Date_of_Journey=data.Date_of_Journey.astype('int64')

#data.Month=odata.Month.astype('int64')

#data.Year=data.Year.astype('int64')


#data.drop(index=6474,inplace=True,axis=0)


data.Travel_Hours=data.Travel_Hours.astype('int64')


categorical=['Airline','Source','Destination','Additional_Info','City1']

numerical=['Total_Stops','Date','Month','Year','Dep_Time_Hour','Dep_Time_Min','Ar
rival_date','Arrival_Time_Hour','Arrival_Time_Mins','Travel_Hours','Travel_Mins']


from sklearn.preprocessing import LabelEncoder

le=LabelEncoder()
```

```python
data.Airline=le.fit_transform(data.Airline)

data.Source=le.fit_transform(data.Source)

data.Destination=le.fit_transform(data.Destination)

data.Total_Stops=le.fit_transform(data.Total_Stops)

data.City1=le.fit_transform(data.City1)

data.City2=le.fit_transform(data.City2)

data.City3=le.fit_transform(data.City3)

data.Addditional_Info=le.fit_transform(data.Additional_Info)

data.head()



import seaborn as sns

c=1

plt.figure(figsize=(20,45))


for i in categorical:

    plt.subplot(6,3,c)

    sns.countplot(data[i])

    plt.xticks(rotation=90)

    plt.tight_layout(pad=3.0)

    c=c+1


plt.show()



#Distribution of 'PRICE' Column
```

```python
plt.figure(figsize=(15,8))

sns.distplot(data.Price)


sns.heatmap(data.corr(annot=True))


#Detecting the Outliers

import seaborn as sns

sns.boxplot(data['Price'])


y=data['Price']

x=data.drop(columns=['Price'],axis=1)


from sklearn.preprocessing import StandardScaler

ss=StandardScaler()


x_scaled=ss.fit_transform


x_scaled=pd.DataFrame(x_scaled,columns=x.columns)

x_scaled.head()


from sklearn.model_selection import train_test_split

x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=42)


x_train.head()


from sklearn.ensemble import RandomForestRegressor,GradientBoostingRegressor,AdaBoostRegressor
```

```python
rfr=RandomForestsRegressor()

gb=GradientBoostingRegressor()

ad=AdaBoostRegressor()


from sklearn.metrics import r2_Score,mean_absolute_error,mean_squared_error


for i in [rfr,gb,ad]:

    i.fit(x_train,y_train)

    y_pred=i.predict(x_test)

    test_score=r2_score(y_test,y_pred)

    train_score=r2_score(y_train,i.predict(x_train))

    if abs(train_score-test_score)<=0.2:

        print(i)


        print("R2 score is",r2_score(y_test,y_pred))

        print("R2 for train data",r2_score(y_train,i.predict(x_train))

        print("Mean Absolue Error is",mean_absolute_error(y_pred,y_test))

        print("Mean Squared Error is",mean_squared_error(y_pred,y_test))

        print("Root Mean Sqaured Error is",(mean_Squared_error(y_pred,y_test,squar
ed-False)))




from sklearn.neighbors import KNeighborsRegressor

from sklearn.svm import SVR

from sklearn.tree import DecisionTreeRegressor
```

```python
from sklearn.metrics import r2_score,mean_absolute_error,mean_squared_error


knn=KNeighborsRegressor()

svr=SVR()

dt=DecisionTreeRegressor()


for i in [knn,svr,dt]:

    i.fit(x_train,y_train)

    y_pred-i.predict(x_test)

    test_score=r2_score(y_test,y_pred)

    train_score=r2score(y_train,i.predict(x_train))

    if abs(train_score-test_score)<=0.1:

        print(i)
print('R2 Score is',r2_score(y_test,y_pred))

print('R2 Score for train data',r2_score(y_train,i,predict(x_train)))

print('Mean Absolute Error is',mean_absolute_error(y_test,y_pred))

print('Mean squared Error is',mean_squared_error(y_test,y_pred))

print('Root Mean Squared Error is',(mean_squared_error(y_test,y_pred,squared=False)))



from sklearn.model_selection import cross_val_score

for i in range(2,5):

    cv=cross_val_score(rfr,x,y,cv=i)

    print(rfr,cv.mean())
```

```python
rfr=RandomForestRegressor(n_estimators=10,max_features='sqrt',max_depth=None)

rfr.fit(x_train,y_train)

y_train_pred=rfr.predict(x_train)

y_test_pred=rfr.predict(x_test)

print("train accurancy",r2_score(y_train_pred,y_train))

print("test accurancy",r2_score(y_test_pred,y_test))


knn=KNeighboursRegressor(n_neighbors=2,algorithm='auto',metric_params=None,n_jobs
=-1)

Knn.fit(x_train,y_train)

y_train_pred=knn.predict (x_train)

y_test_pred=Knn.predict(x_test)

print("train accurancy",r2_score(y_train_pred,y_train))

print("test accurancy",r2_score(y_test_pred,y_test))


import pickle

pickle.dump(rfr,open('model1.pkl','wb'))


from Flask import Flask render_template,'request'

import numpy as np

import pickle


model=pickle.load(open(r"mode[].prl",'rb'))
```

```python
@app,route("/home")

def home():

    return render_template('home.html')


@app.route("/predict")

def home1():

    return render_template('predict.html')


@app.route("/pred",methods=['POST','GET'])

def predict():

    x=[[int(x) for xin request.form.values()]]

    print(x)


    x=np.array(x)

    print(x.shape)


    print(x)

    pred=model.predict(x)

    print(pred)

    return render_template('submit.html',predicition_text=pred)


if __name__ == "__main__":

    app.run(debug=False)


from tensortflow.kers.models import Sequential

from tensortflow.keras.layers import Dense
```

```python
#Fitting the model to the training sets

model = Sequential()


x_train.shape

 (4457, 7163)


model.add(Dense(units =x_train_res.shape[1],activation="relu",kernel_initializer=
"random_uniform"))


model.add(Dense(units=100,activation="relu",kernel_initializer="random_uniform"))


model.add(Dense(units=100,activation="relu",kernel_initializer="random_uniform"))


model.add(Dense(units*1,activation="sigmoid"))


model.compile(optimizer="adam",loss="binary_crossentropy",metrics=['accuracy'])


generator=model.fit(x_train_res,y_train_res,epochs=10,steps_per_epoch=len(x_train
_res)//64)


from sklearn.naive_bayes import MultinomialNB

model=multinomialNB()


#Fitting the model to the training sets

model.fit(x_train_res,y_train_res)
```