```
In [ ]:   1  #*********************************************************************************
          2  #
          3  #                         Session21 Machine Learning 3 Assignment
          4  #
          5  #*********************************************************************************
          6  # Problem Statement : Predict Classification for Each Women whether having affair
          7
          8  # Step 0 : importing the librabries
          9  import numpy as np
         10  import pandas as pd
         11  import statsmodels.api as sm
         12  import seaborn as sns
         13  import matplotlib.pyplot as plt
         14  % matplotlib inline
         15
         16  from patsy import dmatrices
         17  from sklearn.linear_model import LogisticRegression
         18  from sklearn.model_selection import train_test_split ,cross_val_score
         19  from sklearn import metrics
         20  #from sklearn.cross_validation import cross_val_score
         21
         22
         23  # Step1 : Loading the Data
         24
         25  # Description of Variables
         26  # The dataset contains 6366 observations of 9 variables:
         27  # rate_marriage: woman's rating of her marriage (1 = very poor, 5 = very good)
         28  # age: woman's age
         29  # yrs_married: number of years married
         30  # children: number of children
         31  # religious: woman's rating of how religious she is (1 = not religious, 4 = strongly religious)
         32  # educ: level of education (9 = grade school, 12 = high school, 14 = some college, 16 = college graduate, 17 = some g
         33  # occupation: woman's occupation (1 = student, 2 = farming/semi-skilled/unskilled, 3 = "white collar", 4 = teacher/nu
         34  # occupation_husb: husband's occupation (same coding as above)
         35  # affairs: time spent in extra-marital affairs
         36
         37  # Loading the predefined DataSet DAta
         38  dta = sm.datasets.fair.load_pandas().data
         39
         40
         41  print(dta.head())
```

```
42
43  print("\nunique values in Occupation fields \n",'-'*80 )
44  print(dta.occupation.unique())
45  print(dta.occupation_husb.unique())
46
47  dta['affair'] = (dta.affairs > 0).astype(int)
48
49  y, X = dmatrices('affair ~ rate_marriage + age + yrs_married + children + religious + educ + C(occupation) + C(occupa
50                   dta, return_type="dataframe")
51
52
53  print("The number of rows and columns in the Data \n")
54  print(y.shape)
55  print(X.shape)
56  print(X.columns)
57  #print()
58
59  X = X.rename(columns = {'C(occupation)[T.2.0]':'occ_2',
60                          'C(occupation)[T.3.0]':'occ_3',
61                          'C(occupation)[T.4.0]':'occ_4',
62                          'C(occupation)[T.5.0]':'occ_5',
63                          'C(occupation)[T.6.0]':'occ_6',
64                          'C(occupation_husb)[T.2.0]':'occ_husb_2',
65                          'C(occupation_husb)[T.3.0]':'occ_husb_3',
66                          'C(occupation_husb)[T.4.0]':'occ_husb_4',
67                          'C(occupation_husb)[T.5.0]':'occ_husb_5',
68                          'C(occupation_husb)[T.6.0]':'occ_husb_6'})
69  y = np.ravel(y)
70
71  print("The Data with renamed Columns \n",'-'*80)
72  print(X.head(2))
73
74
```

In [ ]:
```python
#Step2 : Analyse the Data

# knowing the size of the data
print(" Analyse the DAta with functions like describe, info , check for nulls \n",'-'*80)
print(X.shape)
print(y.shape)

# knowing the
print(X.describe())
print('-'*80)

print(X.info())

# checking if there are any null fields
print(X.isnull().sum())

print('-'*80,"\n No null values for the Data \n")
```

In [ ]:
```python
# Step3 : Split the Data into Training and Test Set

print(X.columns)
features = ['Intercept','occ_2', 'occ_3', 'occ_4', 'occ_5', 'occ_6', 'occ_husb_2','occ_husb_3', 'occ_husb_4', 'occ_hu
X = X[features]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=1)

print(X_train.shape)
print(X_test.shape)

print(y_train.shape)
print(y_test.shape)

```

In [ ]:
```python
# Step 4 : Doing the Histogram of all variables to visualise the Data better
X_train.hist()

```

In [ ]:
```python
# Step 5: Build the Logistic Regresssion Model and train it

clf_1 = LogisticRegression()
clf_1.fit(X_train, y_train)

# predict the Values
print(clf_1.predict(X_test))


```

In [ ]:
```python
# Step 6 : Model Evaluation
from sklearn.metrics import roc_auc_score, accuracy_score
from sklearn import metrics


accuracy = accuracy_score(y_test, clf_1.predict(X_test))

rocauc= roc_auc_score(y_test, clf_1.predict(X_test))

confusion_matrix = metrics.confusion_matrix(y_test, clf_1.predict(X_test))

classification_report = metrics.classification_report(y_test, clf_1.predict(X_test))

print(accuracy)
print("-"*50)
print(rocauc)
print("-"*50)
print(confusion_matrix)
print("-"*50)
print(classification_report)



```

```
In [ ]:    1  ## Model Evaluation by  cross_val_score
           2
           3  import time
           4  start_time = time.time()
           5  from sklearn.model_selection import cross_val_score
           6  clf_2 = LogisticRegression()
           7
           8  scores= cross_val_score(clf_2, X, y, cv=10)
           9  precision= cross_val_score(clf_2,  X, y, cv=10, scoring='precision')
          10  recall= cross_val_score(clf_2,  X, y, cv=10, scoring='recall')
          11
          12  print(scores.mean())
          13  print(precision.mean())
          14  print(recall.mean())
          15
          16  print("--- %s seconds ---" % (time.time() - start_time))
```

```
In [ ]:    1  # ROC AUC score
           2
           3  from sklearn.metrics import roc_curve, auc
           4  from sklearn import metrics
           5
           6  #X_train
           7  probs = clf_1.predict_proba(X_train)
           8  preds = probs[:,1]
           9
          10
          11  fpr, tpr, threshold = metrics.roc_curve(y_train, preds)
          12  roc_auc = metrics.auc(fpr, tpr)
          13
          14
          15  #X_test
          16  probs1 = clf_1.predict_proba(X_test)
          17  preds1 = probs1[:,1]
          18
          19  fpr1, tpr1, threshold1 = metrics.roc_curve(y_test, preds1)
          20  roc_auc1 = metrics.auc(fpr1, tpr1)
```

In [ ]:

```python
import matplotlib.pyplot as plt

plt.title('Receiver Operating Characteristic')
plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
plt.plot(fpr1, tpr1, 'g', label = 'AUC = %0.2f' % roc_auc1)

plt.legend(['training','test'],loc = 'lower right')

plt.plot([0, 1], [0, 1],'r--')

plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')

plt.show()
```