

In [5]:

```

1  # Problem - I decided to treat this as a classification problem by creating
2  # (did the woman have at least one affair?) and trying to predict the classi
3  # Dataset
4  # The dataset I chose is the affairs dataset that comes with Statsmodels. It
5  # by Redbook magazine, in which married women were asked about their partici
6  # information about the study is available in a 1978 paper from the Journal
7  # Description of Variables
8  # The dataset contains 6366 observations of 9 variables:
9  # rate_marriage: woman's rating of her marriage (1 = very poor, 5 = very goo
10 # age: woman's age
11 # yrs_married: number of years married
12 # children: number of children
13 # religious: woman's rating of how religious she is (1 = not religious, 4 =
14 # educ: level of education (9 = grade school, 12 = high school, 14 = some co
15 # college graduate, 17 = some graduate school, 20 = advanced degree)
16 # occupation: woman's occupation (1 = student, 2 = farming/semi-skilled/unsk
17 # "white collar", 4 = teacher/nurse/writer/technician/skilled, 5 = manageria
18 # professional with advanced degree)
19 # occupation_husb: husband's occupation (same coding as above)
20 # affairs: time spent in extra-marital affairs
21 # Code to Loading data and modules
22 import numpy as np
23 import pandas as pd
24 import statsmodels.api as sm
25 import matplotlib.pyplot as plt
26 from patsy import dmatrices
27 from sklearn.linear_model import LogisticRegression
28 from sklearn.cross_validation import train_test_split
29 from sklearn import metrics
30 from sklearn.cross_validation import cross_val_score
31 dta = sm.datasets.fair.load_pandas().data
32 # add "affair" column: 1 represents having affairs, 0 represents not
33 dta['affair'] = (dta.affairs > 0).astype(int)
34 y, X = dmatrices('affair ~ rate_marriage + age + yrs_married + children + \
35 religious + educ + C(occupation) + C(occupation_husb)',
36 dta, return_type="dataframe")
37 X = X.rename(columns = {'C(occupation)[T.2.0]': 'occ_2',
38 'C(occupation)[T.3.0]': 'occ_3',
39 'C(occupation)[T.4.0]': 'occ_4',
40 'C(occupation)[T.5.0]': 'occ_5',
41 'C(occupation)[T.6.0]': 'occ_6',
42 'C(occupation_husb)[T.2.0]': 'occ_husb_2',
43 'C(occupation_husb)[T.3.0]': 'occ_husb_3',
44 'C(occupation_husb)[T.4.0]': 'occ_husb_4',
45 'C(occupation_husb)[T.5.0]': 'occ_husb_5',
46 'C(occupation_husb)[T.6.0]': 'occ_husb_6'})
47 y = np.ravel(y)

```

```
In [4]: 1 import numpy as np
2 import pandas as pd
3 #using pandas.tseries instead of statsmodels.api
4 import pandas.tseries as pdt
5 import matplotlib.pyplot as plt
6 from patsy import dmatrices
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.model_selection import train_test_split
9 from sklearn import metrics
10 from sklearn.model_selection import cross_val_score
11 #To avoid warnings
12 import warnings
13 warnings.filterwarnings('ignore')
14 dta = sm.datasets.fair.load_pandas().data
15 df_affair = dta.copy()
```

```
In [6]: 1 # add "affair" column: 1 represents having affairs, 0 represents not
2
3 dta['affair'] = (dta.affairs > 0).astype(int)
4 y, X = dmatrices('affair ~ rate_marriage + age + yrs_married + children + \
5 religious + educ + C(occupation) + C(occupation_husb)',
6 dta, return_type="dataframe")
7
8 X = X.rename(columns = {'C(occupation)[T.2.0]': 'occ_2',
9 'C(occupation)[T.3.0]': 'occ_3',
10 'C(occupation)[T.4.0]': 'occ_4',
11 'C(occupation)[T.5.0]': 'occ_5',
12 'C(occupation)[T.6.0]': 'occ_6',
13 'C(occupation_husb)[T.2.0]': 'occ_husb_2',
14 'C(occupation_husb)[T.3.0]': 'occ_husb_3',
15 'C(occupation_husb)[T.4.0]': 'occ_husb_4',
16 'C(occupation_husb)[T.5.0]': 'occ_husb_5',
17 'C(occupation_husb)[T.6.0]': 'occ_husb_6'})
18 y = np.ravel(y)
```

```
In [7]: 1 dta.head()
```

Out[7]:

	rate_marriage	age	yrs_married	children	religious	educ	occupation	occupation_husb	affa
0	3.0	32.0	9.0	3.0	3.0	17.0	2.0	5.0	0.111
1	3.0	27.0	13.0	3.0	1.0	14.0	3.0	4.0	3.2307
2	4.0	22.0	2.5	0.0	1.0	16.0	3.0	5.0	1.4000
3	4.0	37.0	16.5	4.0	3.0	16.0	5.0	5.0	0.7272
4	5.0	27.0	9.0	1.0	1.0	14.0	3.0	4.0	4.6666

```
In [8]: 1 dta.shape
```

Out[8]: (6366, 10)

In [9]: 1 X.head()

Out[9]:

	Intercept	occ_2	occ_3	occ_4	occ_5	occ_6	occ_husb_2	occ_husb_3	occ_husb_4	occ_husb_5
0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	1.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
2	1.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	1.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0
4	1.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0

In [10]: 1 y

Out[10]: array([1., 1., 1., ..., 0., 0., 0.])

In [11]: 1 print("Lets analyze the data and look at the summary statistics")
2 dta.describe()

Lets analyze the data and look at the summary statistics

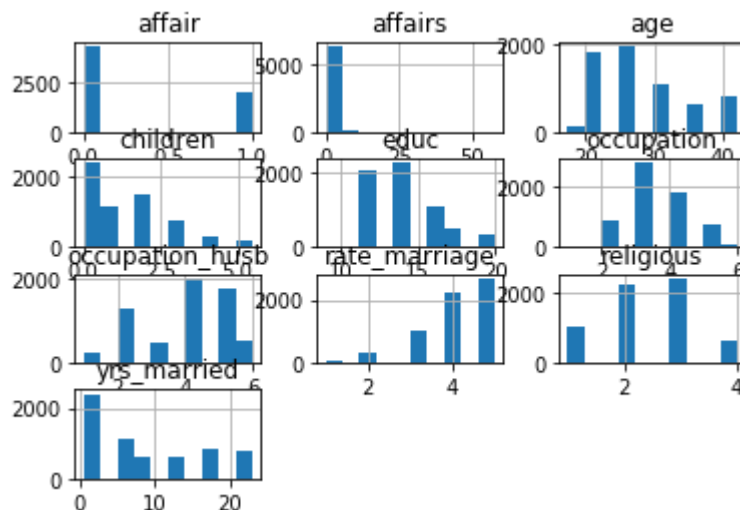
Out[11]:

	rate_marriage	age	yrs_married	children	religious	educ	occupatio
count	6366.000000	6366.000000	6366.000000	6366.000000	6366.000000	6366.000000	6366.000000
mean	4.109645	29.082862	9.009425	1.396874	2.426170	14.209865	3.42412
std	0.961430	6.847882	7.280120	1.433471	0.878369	2.178003	0.94239
min	1.000000	17.500000	0.500000	0.000000	1.000000	9.000000	1.00000
25%	4.000000	22.000000	2.500000	0.000000	2.000000	12.000000	3.00000
50%	4.000000	27.000000	6.000000	1.000000	2.000000	14.000000	3.00000
75%	5.000000	32.000000	16.500000	2.000000	3.000000	16.000000	4.00000
max	5.000000	42.000000	23.000000	5.500000	4.000000	20.000000	6.00000

```
In [12]: 1 # plot all of the columns
2 %matplotlib inline
3 plt.figure(figsize=(20,18))
4 dta.hist()
```

```
Out[12]: array([[<matplotlib.axes._subplots.AxesSubplot object at 0x09FBAE90>,
<matplotlib.axes._subplots.AxesSubplot object at 0x09F56D50>,
<matplotlib.axes._subplots.AxesSubplot object at 0x09F46C50>],
[<matplotlib.axes._subplots.AxesSubplot object at 0x09F3AB50>,
<matplotlib.axes._subplots.AxesSubplot object at 0x09F60A50>,
<matplotlib.axes._subplots.AxesSubplot object at 0x09F689F0>],
[<matplotlib.axes._subplots.AxesSubplot object at 0x0A02CB30>,
<matplotlib.axes._subplots.AxesSubplot object at 0x0A0465F0>,
<matplotlib.axes._subplots.AxesSubplot object at 0x0A046B70>],
[<matplotlib.axes._subplots.AxesSubplot object at 0x0A05FCF0>,
<matplotlib.axes._subplots.AxesSubplot object at 0x039C5DB0>,
<matplotlib.axes._subplots.AxesSubplot object at 0x039E0E70>]],
dtype=object)
```

<Figure size 1440x1296 with 0 Axes>



```
In [13]: 1 print("Split the data into training and test set")
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, ran
4 print(X_train.shape)
5 print(y_train.shape)
6 print(X_test.shape)
7 print(y_test.shape)
```

```
Split the data into training and test set
(4456, 17)
(4456,)
(1910, 17)
(1910,)
```

```
In [14]: 1 #We will use the statsmodels Logit function for logistic regression
          2 logit = sm.Logit(y_train, X_train)
          3
          4 # fit the model
          5 result = logit.fit()
```

```
Optimization terminated successfully.
      Current function value: 0.544479
      Iterations 6
```

```
In [15]: 1 predictions = result.predict(X_test)
          2 predictions
```

```
Out[15]: 2764    0.653211
          4481    0.087718
          5360    0.273074
          5802    0.249471
          1220    0.249630
          5812    0.166215
          3719    0.160619
          3848    0.202858
          1865    0.760648
          2535    0.310242
          2505    0.104535
          6273    0.186280
          3710    0.075798
          4229    0.300914
          1262    0.736723
          5321    0.593884
          3790    0.296514
          994     0.732196
          5644    0.296810
          2252    0.156072
          1804    0.203707
          861     0.448163
          1601    0.089842
          1718    0.471631
          2976    0.168971
          3603    0.161252
          4130    0.396058
          5824    0.361373
          5901    0.252252
          4408    0.087566
          ...
          5615    0.242151
          1737    0.515613
          2701    0.185493
          4024    0.419446
          1012    0.129569
          3888    0.143331
          4746    0.207049
          5607    0.132057
          1946    0.874461
          3119    0.133824
          156     0.604838
          1752    0.653266
          624     0.612704
          4622    0.556290
          1788    0.693977
          500     0.264134
          726     0.220550
          4162    0.087718
          48      0.158934
          1691    0.572606
          5882    0.152058
          2244    0.395267
          1985    0.841614
```

```

2853    0.223091
18      0.803795
3053    0.144139
1875    0.207506
5851    0.437646
4962    0.190124
1995    0.249630
Length: 1910, dtype: float64

```

```

In [16]: 1 from scipy import stats
         2 stats.chisqprob = lambda chisq, df: stats.chi2.sf(chisq, df)

```

```

In [17]: 1 result.summary()

```

Out[17]:

Logit Regression Results

Dep. Variable:	y	No. Observations:	4456
Model:	Logit	Df Residuals:	4439
Method:	MLE	Df Model:	16
Date:	Tue, 05 Mar 2019	Pseudo R-squ.:	0.1360
Time:	11:51:13	Log-Likelihood:	-2426.2
converged:	True	LL-Null:	-2808.3
		LLR p-value:	2.844e-152

	coef	std err	z	P> z	[0.025	0.975]
Intercept	2.4842	0.777	3.198	0.001	0.961	4.007
occ_2	0.9414	0.658	1.432	0.152	-0.347	2.230
occ_3	1.2324	0.652	1.890	0.059	-0.046	2.511
occ_4	0.9731	0.653	1.490	0.136	-0.307	2.254
occ_5	1.6017	0.657	2.436	0.015	0.313	2.890
occ_6	1.8242	0.707	2.581	0.010	0.439	3.209
occ_husb_2	0.0649	0.215	0.302	0.762	-0.356	0.486
occ_husb_3	0.1976	0.235	0.841	0.400	-0.263	0.658
occ_husb_4	0.0304	0.208	0.146	0.884	-0.377	0.438
occ_husb_5	-0.0052	0.210	-0.025	0.980	-0.417	0.406
occ_husb_6	-0.0183	0.236	-0.078	0.938	-0.481	0.445
rate_marriage	-0.7145	0.038	-18.929	0.000	-0.788	-0.640
age	-0.0577	0.012	-4.686	0.000	-0.082	-0.034
yrs_married	0.1081	0.013	8.243	0.000	0.082	0.134
children	-0.0126	0.038	-0.329	0.742	-0.088	0.062
religious	-0.3889	0.042	-9.342	0.000	-0.470	-0.307
educ	0.0046	0.021	0.224	0.823	-0.036	0.045

```
In [18]: 1 print("Logistic Regression with scikit-learn")
          2 dta.head()
```

Logistic Regression with scikit-learn

Out[18]:

	rate_marriage	age	yrs_married	children	religious	educ	occupation	occupation_husb	affair
0	3.0	32.0	9.0	3.0	3.0	17.0	2.0	5.0	0.111
1	3.0	27.0	13.0	3.0	1.0	14.0	3.0	4.0	3.2307
2	4.0	22.0	2.5	0.0	1.0	16.0	3.0	5.0	1.4000
3	4.0	37.0	16.5	4.0	3.0	16.0	5.0	5.0	0.7273
4	5.0	27.0	9.0	1.0	1.0	14.0	3.0	4.0	4.6667

```
In [19]: 1 print('Exploratory data analysis')
          2 # people having affair is represented with 1 and not having affair is repres
          3 dta.affair.value_counts()
```

Exploratory data analysis

Out[19]: 0 4313
1 2053
Name: affair, dtype: int64

```
In [20]: 1 print("We can conclude that women who have affairs, rate their marriage lower")
          2 dta.groupby('affair').mean()
```

We can conclude that women who have affairs, rate their marriage lower based on our findings from below table

Out[20]:

	rate_marriage	age	yrs_married	children	religious	educ	occupation	occupation_husb	affair
0	4.329701	28.390679	7.989335	1.238813	2.504521	14.322977	3.405286	3.405286	0
1	3.647345	30.537019	11.152460	1.728933	2.261568	13.972236	3.463712	3.463712	1


```
In [21]: 1 print('Checking rate_marriage parameter')
          2 print('We can say with an increase in age, yrs_married and children correlate')
          3 dta.groupby('rate_marriage').mean()
```

Checking rate_marriage parameter

We can say with an increase in age, yrs_married and children correlate with increase in affairs based on findings.

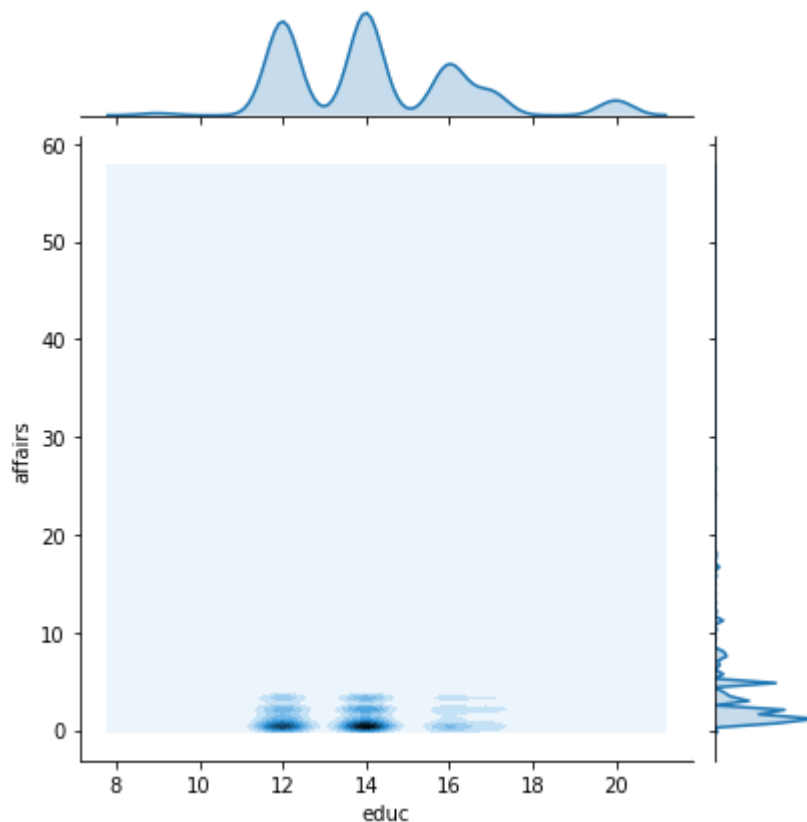
Out[21]:

		age	yrs_married	children	religious	educ	occupation	occupation_husb
rate_marriage								
	1.0	33.823232	13.914141	2.308081	2.343434	13.848485	3.232323	3.838384
	2.0	30.471264	10.727011	1.735632	2.330460	13.864943	3.327586	3.764368
	3.0	30.008056	10.239174	1.638469	2.308157	14.001007	3.402820	3.798590
	4.0	28.856601	8.816905	1.369536	2.400981	14.144514	3.420161	3.835861
	5.0	28.574702	8.311662	1.252794	2.506334	14.399776	3.454918	3.892697

```
In [22]: 1 print('Lets visualize our data')
          2 import seaborn as sns
          3 sns.jointplot(x='educ',y='affairs',data=dta,kind='kde')
```

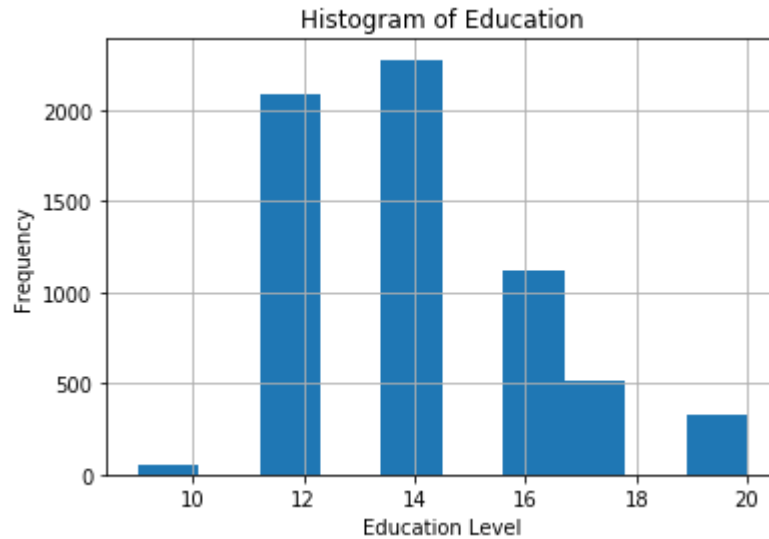
Lets visualize our data

Out[22]: <seaborn.axisgrid.JointGrid at 0x9eae170>



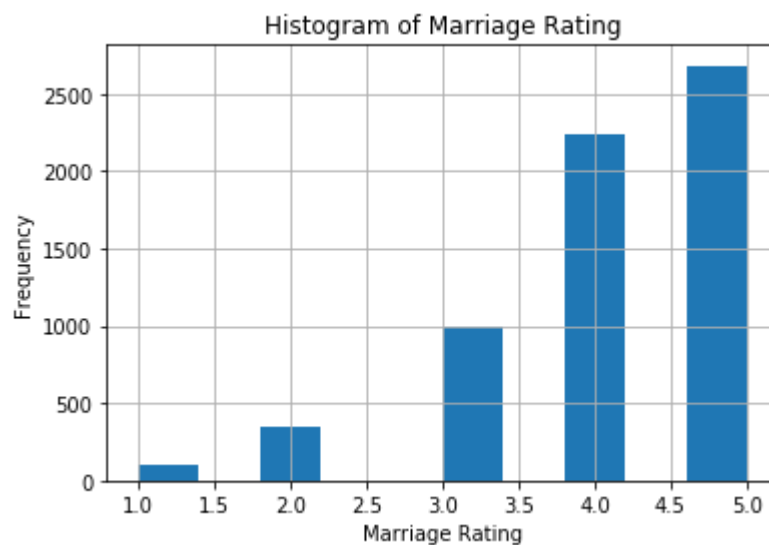
```
In [23]: 1 # histogram of education
2 dta.educ.hist()
3 plt.title('Histogram of Education')
4 plt.xlabel('Education Level')
5 plt.ylabel('Frequency')
```

Out[23]: Text(0, 0.5, 'Frequency')



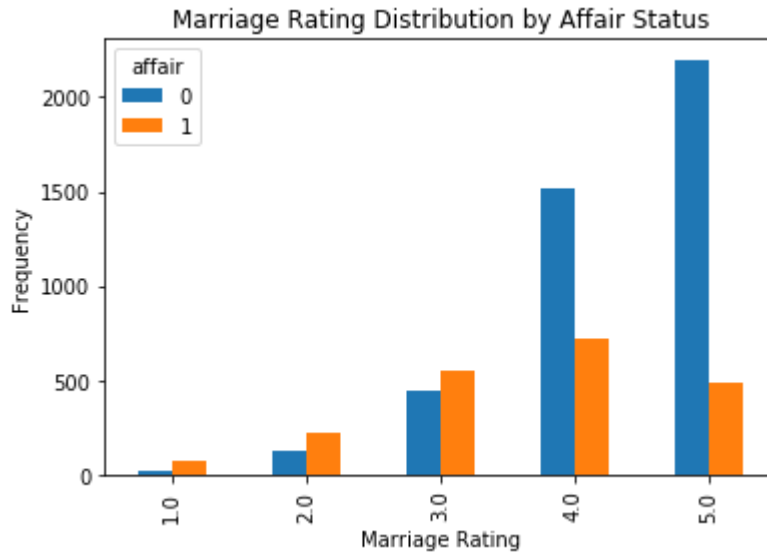
```
In [24]: 1 # histogram of marriage rating
2 dta.rate_marriage.hist()
3 plt.title('Histogram of Marriage Rating')
4 plt.xlabel('Marriage Rating')
5 plt.ylabel('Frequency')
```

Out[24]: Text(0, 0.5, 'Frequency')



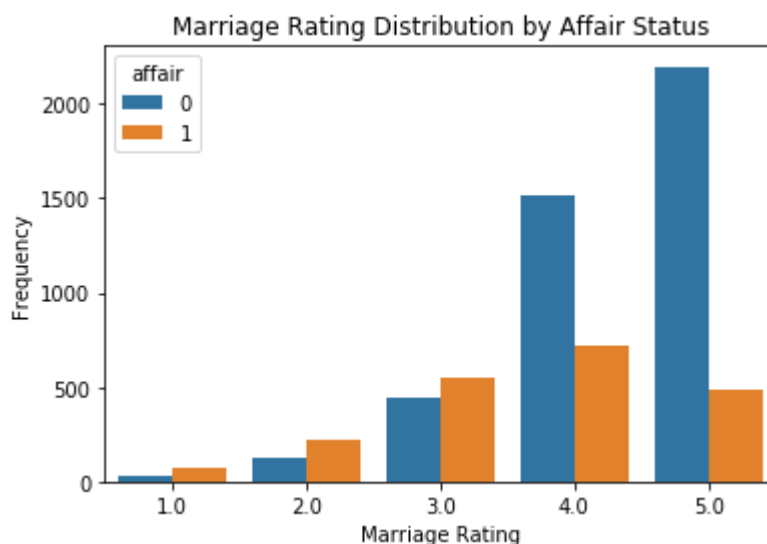
```
In [25]: 1 # barplot of marriage rating grouped by affair (True or False)
2 pd.crosstab(dta.rate_marriage, dta.affair).plot(kind='bar')
3 plt.title('Marriage Rating Distribution by Affair Status')
4 plt.xlabel('Marriage Rating')
5 plt.ylabel('Frequency')
```

Out[25]: Text(0, 0.5, 'Frequency')



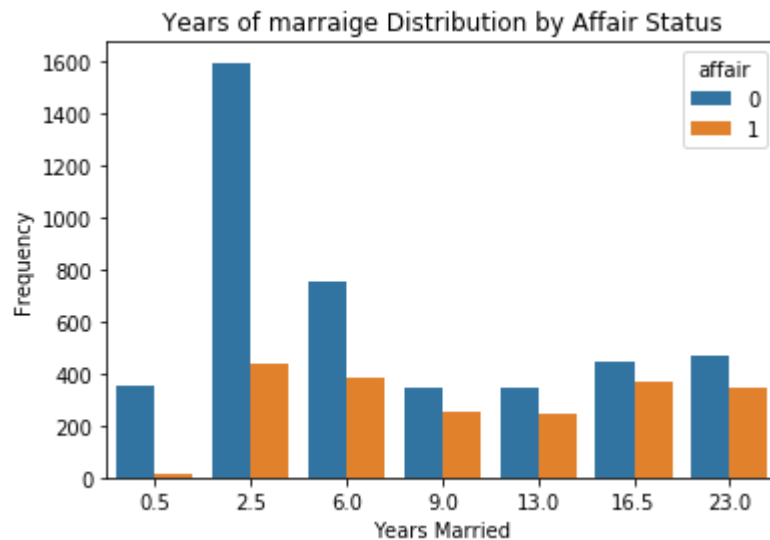
```
In [26]: 1 sns.countplot(x='rate_marriage',data=dta,hue='affair')
2 plt.title('Marriage Rating Distribution by Affair Status')
3 plt.xlabel('Marriage Rating')
4 plt.ylabel('Frequency')
```

Out[26]: Text(0, 0.5, 'Frequency')



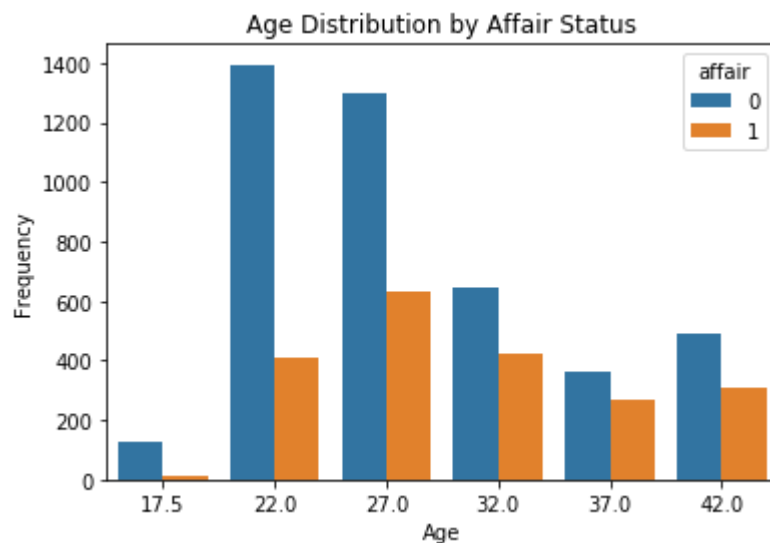
```
In [27]: 1 sns.countplot(x='yrs_married',data=dta,hue='affair')
2 plt.title('Years of marriage Distribution by Affair Status')
3 plt.xlabel('Years Married')
4 plt.ylabel('Frequency')
```

Out[27]: Text(0, 0.5, 'Frequency')



```
In [28]: 1 import seaborn as sns
2 sns.countplot(x='age',data=dta,hue='affair')
3 plt.title('Age Distribution by Affair Status')
4 plt.xlabel('Age')
5 plt.ylabel('Frequency')
```

Out[28]: Text(0, 0.5, 'Frequency')



```
In [29]: 1 print("Model Evaluation Using a Validation Set")
2 from sklearn.model_selection import train_test_split
3 # evaluate the model by splitting into train and test sets
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, ran
5 print(X_train.shape)
6 print(y_train.shape)
7 print(X_test.shape)
8 print(y_test.shape)
```

Model Evaluation Using a Validation Set

(4456, 17)

(4456,)

(1910, 17)

(1910,)

```
In [30]: 1 model = LogisticRegression()
2 model.fit(X_train, y_train)
```

```
Out[30]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=100, multi_class='warn',
n_jobs=None, penalty='l2', random_state=None, solver='warn',
tol=0.0001, verbose=0, warm_start=False)
```

```
In [31]: 1 print(model.score(X_train,y_train))
2 print("Training set has 73% accuracy")
```

0.723967684021544

Training set has 73% accuracy

```
In [32]: 1 print("Use the test data set to predict the class / labels")
2 # predict class labels for the test set
3 predicted = model.predict(X_test)
4 predicted
```

Use the test data set to predict the class / labels

```
Out[32]: array([1., 0., 0., ..., 0., 0., 0.])
```

```
In [33]: 1 # generate class probabilities
2 probs = model.predict_proba(X_test)
3 probs
```

```
Out[33]: array([[0.35146338, 0.64853662],
[0.90955084, 0.09044916],
[0.72567333, 0.27432667],
...,
[0.55727384, 0.44272616],
[0.81207046, 0.18792954],
[0.74734603, 0.25265397]])
```

```
In [34]: 1 print('Evaluating the model')
2 # generate evaluation metrics
3 print(metrics.accuracy_score(y_test,predicted))
4 print(metrics.roc_auc_score(y_test, probs[:, 1]))
5 print("The accuracy of the model is 73% similar to the training data.")
```

Evaluating the model

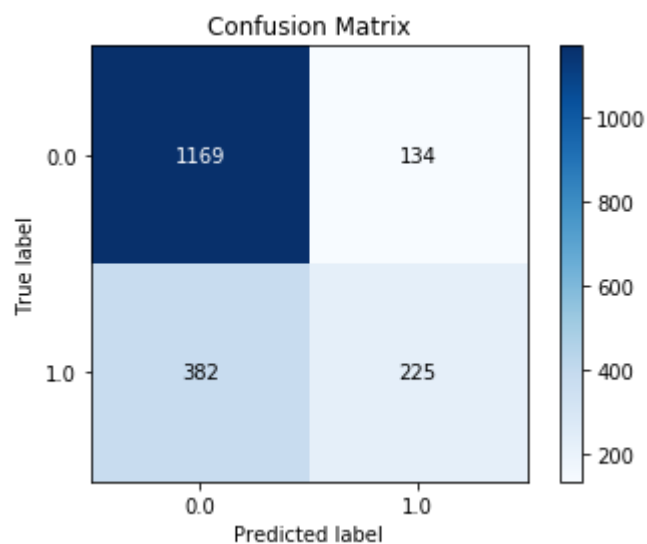
0.7298429319371728

0.745950606950631

The accuracy of the model is 73% similar to the training data.

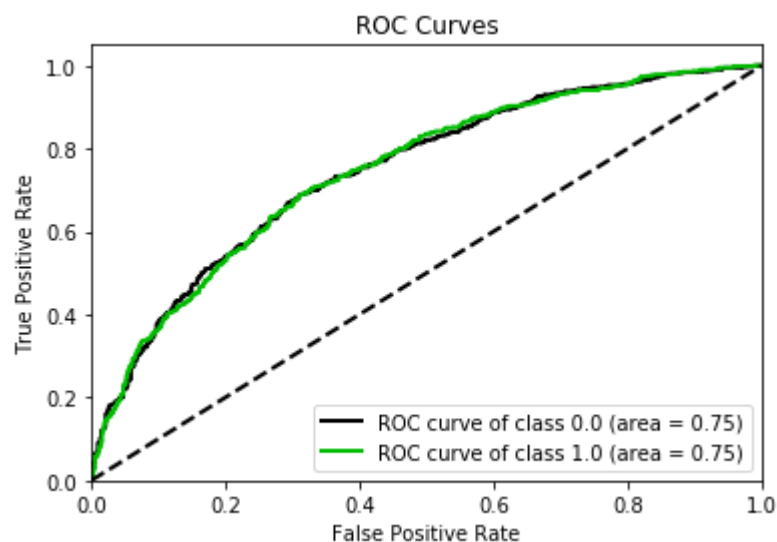
```
In [35]: 1 #Using confusion matrix to describe the performance of the classification mo
2 import scikitplot
3 scikitplot.metrics.plot_confusion_matrix(y_test,predicted)
```

Out[35]: <matplotlib.axes._subplots.AxesSubplot at 0xa8acf50>



```
In [36]: 1 # Plotting the true positive rate (TPR) against the false positive rate (FPR)
2 scikitplot.metrics.plot_roc_curve(y_test, probs,curves=['each_class'])
```

Out[36]: <matplotlib.axes._subplots.AxesSubplot at 0xa9a1dd0>



```
In [39]: 1 #accuracy report
2 print(metrics.classification_report(y_test, predicted))
```

	precision	recall	f1-score	support	
	0.0	0.75	0.90	0.82	1303
	1.0	0.63	0.37	0.47	607
micro avg	0.73	0.73	0.73		1910
macro avg	0.69	0.63	0.64		1910
weighted avg	0.71	0.73	0.71		1910

```
In [40]: 1 from sklearn.metrics import confusion_matrix
2 cf = confusion_matrix(y_test,predicted)
3 type(cf)
```

Out[40]: numpy.ndarray

```
In [41]: 1 cf.shape
```

Out[41]: (2, 2)

```
In [43]: 1 #Calculation of Precision Recall and F1 score
2 TN = cf[0,0] #True Negative
3 FP = cf[0,1] #False Positive
4 FN = cf[1,0] #False Negative
5 TP = cf[1,1] #True Positive
6
7 Precision = TP / (TP + FP)
8 Recall = TP / (TP + FN)
9 F1 = (2 *(Precision * Recall)) / (Precision + Recall)
10 print("Precision : {} , Recall : {}, F1 : {}".format(Precision,Recall,F1))
11
```

Precision : 0.6267409470752089 , Recall : 0.37067545304777594, F1 : 0.4658385093167702

```
In [44]: 1 #Calculation of True Positive Rate and False Positive Rate
2 TPR = (TP) / (TP + FN ) #equal to Recall
3 FPR = FP / (FP + TN )
4 print("True Positive Rate : {}, False Positive Rate : {}".format(TPR,FPR))
5
```

True Positive Rate : 0.37067545304777594, False Positive Rate : 0.10283960092095165

```
In [45]: 1 # evaluate the model using 10-fold cross-validation
2 scores = cross_val_score(LogisticRegression(), X, y, scoring='accuracy', cv=
3 scores, scores.mean())
```

Out[45]: (array([0.72100313, 0.70219436, 0.73824451, 0.70597484, 0.70597484,
0.72955975, 0.7327044 , 0.70440252, 0.75157233, 0.75
0.7241630685514876])

```
In [46]: 1 print('Predicting the Probability of an Affair')
          2 print(model.predict_proba(np.array([[1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 3, 25,
          3 print('The predicted probability of an affair is 23%')
```

Predicting the Probability of an Affair

[[0.77301481 0.22698519]]

The predicted probability of an affair is 23%

```
In [48]: 1 # Let's predict the probability of an affair for a random woman not present
          2 # She's a 30-year-old teacher who graduated college, has been married for 10
          3 # as strongly religious, rates her marriage as fair, and her husband is a fa
          4 print(model.predict_proba(np.array([[1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 3, 30,
          5 print('The predicted probability of an affair is 31%')
```

[[0.68617099 0.31382901]]

The predicted probability of an affair is 31%