

```
import itertools
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.ticker import NullFormatter
import pandas as pd
import numpy as np
import matplotlib.ticker as ticker
from sklearn import preprocessing
%matplotlib inline
```

```
!wget -O loan_train.csv https://s3-api.us-gEO.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_train.csv

--2023-01-22 11:21:47-- https://s3-api.us-gEO.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_train.csv
Resolving s3-api.us-gEO.objectstorage.softlayer.net (s3-api.us-gEO.objectstorage.softlayer.net)... 67.228.254.196
Connecting to s3-api.us-gEO.objectstorage.softlayer.net (s3-api.us-gEO.objectstorage.softlayer.net)|67.228.254.196|:443... connecte
HTTP request sent, awaiting response... 200 OK
Length: 23101 (23K) [text/csv]
Saving to: 'loan_train.csv'

loan_train.csv      100%[=====>]  22.56K  --.-KB/s    in 0.02s

2023-01-22 11:21:47 (1.04 MB/s) - 'loan_train.csv' saved [23101/23101]
```

```
df = pd.read_csv('loan_train.csv')
df.head(10)
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	0	0	PAIDOFF	1000	30	9/8/2016	10/7/2016	45	High School or Below	male
1	2	2	PAIDOFF	1000	30	9/8/2016	10/7/2016	33	Bechalar	female
2	3	3	PAIDOFF	1000	15	9/8/2016	9/22/2016	27	college	male
3	4	4	PAIDOFF	1000	30	9/9/2016	10/8/2016	28	college	female
4	6	6	PAIDOFF	1000	30	9/9/2016	10/8/2016	29	college	male
5	7	7	PAIDOFF	1000	30	9/9/2016	10/8/2016	36	college	male
6	8	8	PAIDOFF	1000	30	9/9/2016	10/8/2016	28	college	male
7	9	9	PAIDOFF	800	15	9/10/2016	9/24/2016	26	college	male
8	10	10	PAIDOFF	300	7	9/10/2016	9/16/2016	29	college	male

```
df.shape
```

(346, 10)

```
df['due_date'] = pd.to_datetime(df['due_date'])
df['effective_date'] = pd.to_datetime(df['effective_date'])
df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	High School or Below	male
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bechalar	female
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	male

```
df['loan_status'].value_counts()
```

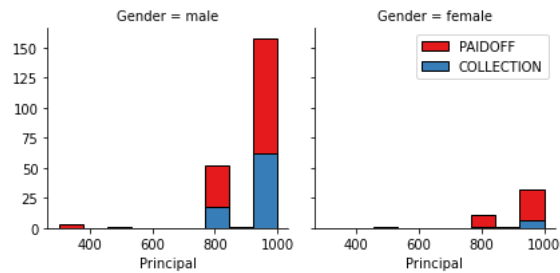
```
PAIDOFF      260
COLLECTION    86
Name: loan_status, dtype: int64
```

```
import seaborn as sns

bins = np.linspace(df.Principal.min(), df.Principal.max(), 10)
```

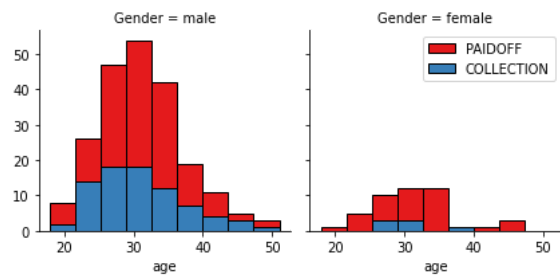
```
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)
g.map(plt.hist, 'Principal', bins=bins, ec="k")

g.axes[-1].legend()
plt.show()
```

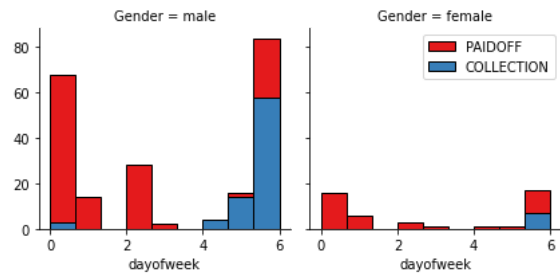


```
bins = np.linspace(df.age.min(), df.age.max(), 10)
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)
g.map(plt.hist, 'age', bins=bins, ec="k")

g.axes[-1].legend()
plt.show()
```



```
df['dayofweek'] = df['effective_date'].dt.dayofweek
bins = np.linspace(df.dayofweek.min(), df.dayofweek.max(), 10)
g = sns.FacetGrid(df, col="Gender", hue="loan_status", palette="Set1", col_wrap=2)
g.map(plt.hist, 'dayofweek', bins=bins, ec="k")
g.axes[-1].legend()
plt.show()
```



```
df['weekend'] = df['dayofweek'].apply(lambda x: 1 if (x>3) else 0)
df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender	dayofweek
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	High School or Below	male	3
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bechalar	female	3
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	male	3

```
df.groupby(['Gender'])['loan_status'].value_counts(normalize=True)

Gender  loan_status
female  PAIDOFF      0.865385
        COLLECTION   0.134615
male    PAIDOFF      0.731293
        COLLECTION   0.268707
Name: loan_status, dtype: float64
```

```
df['Gender'].replace(to_replace=['male','female'], value=[0,1],inplace=True)
df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender	dayofweek
0	0	0	PAIDOFF	1000	30	2016-09-08	2016-10-07	45	High School or Below	0	3
1	2	2	PAIDOFF	1000	30	2016-09-08	2016-10-07	33	Bechalar	1	3
2	3	3	PAIDOFF	1000	15	2016-09-08	2016-09-22	27	college	0	3

```
df.groupby(['education'])['loan_status'].value_counts(normalize=True)
```

education	loan_status	
Bechalar	PAIDOFF	0.750000
	COLLECTION	0.250000
High School or Below	PAIDOFF	0.741722
	COLLECTION	0.258278
Master or Above	COLLECTION	0.500000
	PAIDOFF	0.500000
college	PAIDOFF	0.765101
	COLLECTION	0.234899

Name: loan_status, dtype: float64

```
df[['Principal','terms','age','Gender','education']].head()
```

	Principal	terms	age	Gender	education
0	1000	30	45	0	High School or Below
1	1000	30	33	1	Bechalar
2	1000	15	27	0	college
3	1000	30	28	1	college
4	1000	30	29	0	college

```
Feature = df[['Principal','terms','age','Gender','weekend']]
Feature = pd.concat([Feature,pd.get_dummies(df['education'])], axis=1)
Feature.drop(['Master or Above'], axis = 1,inplace=True)
Feature.head()
```

	Principal	terms	age	Gender	weekend	Bechalar	High School or Below	college
0	1000	30	45	0	0	0	1	0
1	1000	30	33	1	0	1	0	0
2	1000	15	27	0	0	0	0	1
3	1000	30	28	1	1	0	0	1
4	1000	30	29	0	1	0	0	1

```
X = Feature
X[0:5]
```

	Principal	terms	age	Gender	weekend	Bechalar	High School or Below	college
0	1000	30	45	0	0	0	1	0
1	1000	30	33	1	0	1	0	0
2	1000	15	27	0	0	0	0	1
3	1000	30	28	1	1	0	0	1
4	1000	30	29	0	1	0	0	1

```
df['loan_status'].replace(to_replace=['PAIDOFF','COLLECTION'], value=[0,1],inplace=True)
df.head()
```

```

Unnamed: 0      Unnamed: 0.1  loan_status  Principal  terms  effective_date  due_date  age  education  Gender  dayofweek

```

```

y = df['loan_status'].values
y[0:20]

```

```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
```

```

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split( X, y, test_size=0.40, random_state=5)

```

```

print ('Train set:', X_train.shape, y_train.shape)
print ('Test set:', X_test.shape, y_test.shape)

```

```

Train set: (207, 8) (207,)
Test set: (139, 8) (139,)

```

```

display(X_train.shape, y_train.shape)
display(X_test.shape, y_test.shape)

```

```

↳ (207, 8)
(207,)
(139, 8)
(139,)

```

```

# (technically should be done after train test split )
X_train = preprocessing.StandardScaler().fit(X_train).transform(X_train)
X_test = preprocessing.StandardScaler().fit(X_test).transform(X_test)
X_train[0:5]

```

```

array([[ 0.52533066, -0.97207061,  0.53495582, -0.45109685,  0.85993942,
         2.76134025, -0.86846836, -0.91206272],
       [ 0.52533066,  0.91317564, -0.47326137,  2.21681883,  0.85993942,
        -0.36214298, -0.86846836,  1.09641582],
       [ 0.52533066,  0.91317564, -0.13718897, -0.45109685,  0.85993942,
        -0.36214298,  1.15145243, -0.91206272],
       [ 0.52533066,  0.91317564, -0.64129757, -0.45109685,  0.85993942,
        -0.36214298, -0.86846836,  1.09641582],
       [ 0.52533066, -0.97207061, -0.47326137,  2.21681883, -1.16287262,
        -0.36214298, -0.86846836,  1.09641582]])

```

```

from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score

```

```

# We will be checking for value of K from 0 to 25 so
n = 25
accuracy = np.zeros(n)
for i in range(1,n+1):
    clf = KNeighborsClassifier(n_neighbors = i).fit(X_train, y_train)
    y_test_predicted = clf.predict(X_test)
    accuracy[i-1] = (accuracy_score(y_test, y_test_predicted))
accuracy

```

```

array([0.65467626, 0.68345324, 0.6618705 , 0.71223022, 0.69064748,
       0.69064748, 0.65467626, 0.70503597, 0.64028777, 0.69784173,
       0.69784173, 0.69784173, 0.69064748, 0.69784173, 0.65467626,
       0.69064748, 0.6618705 , 0.67625899, 0.67625899, 0.67625899,
       0.70503597, 0.70503597, 0.71942446, 0.74100719, 0.74820144])

```

```

plt.plot(range(1,n+1),accuracy, 'g')
plt.ylabel('Accuracy')
plt.xlabel('Number of Neighbors (K)')
plt.show()

```

```

accuracy = pd.DataFrame(accuracy)
print("Maximum Accuracy Got is - " )
accuracy.sort_values(by = 0, ascending = False)[0:1]

```



```
clf_KNN = KNeighborsClassifier(n_neighbors = 24).fit(X_train, y_train)
# y_test_pred_KNN = clf_KNN.predict(X_test)
```

```
from sklearn.tree import DecisionTreeClassifier
```

```
clf2 = DecisionTreeClassifier(criterion = 'gini').fit(X_train, y_train)
y_test_pred_KNN = clf.predict(X_test)
print("Accuracy using criterion as gini - ", accuracy_score(y_test, y_test_pred_KNN))
clf3 = DecisionTreeClassifier(criterion = 'entropy').fit(X_train, y_train)
y_test_pred_KNN = clf2.predict(X_test)
print("Accuracy using criterion as entropy - ", accuracy_score(y_test, y_test_pred_KNN))
```

```
Accuracy using criterion as gini - 0.7482014388489209
Accuracy using criterion as entropy - 0.6474820143884892
```

```
# using criterion as gini
clf_DT = DecisionTreeClassifier(criterion = 'gini').fit(X_train, y_train)
```

```
from sklearn.svm import SVC
```

```
clf4 = SVC(kernel = 'poly').fit(X_train, y_train)
print("accuracy using polynomial kernel - ", accuracy_score(y_test, clf2.predict(X_test)))
clf5 = SVC(kernel = 'rbf').fit(X_train, y_train)
print("accuracy using Radial Basis function kernel - ", accuracy_score(y_test, clf3.predict(X_test)))
```

```
accuracy using polynomial kernel - 0.6474820143884892
accuracy using Radial Basis function kernel - 0.6474820143884892
```

```
# using linear for kernel in our SVM model
clf_SVM = SVC(kernel = 'poly', random_state = 4).fit(X_train, y_train)
```

```
from sklearn.linear_model import LogisticRegression
```

```
clf_LR = LogisticRegression(solver='lbfgs', warm_start = True)
clf_LR.fit(X_train, y_train)
#print("accuracy score - ", accuracy_score(y_test, clf_LR.predict(X_test)))
```

```
LogisticRegression(warm_start=True)
```

```
!wget -O loan_test.csv https://s3-api.us-gio.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_test.csv
```

```
--2023-01-22 11:32:53-- https://s3-api.us-gio.objectstorage.softlayer.net/cf-courses-data/CognitiveClass/ML0101ENV3/labs/loan_test
Resolving s3-api.us-gio.objectstorage.softlayer.net (s3-api.us-gio.objectstorage.softlayer.net)... 67.228.254.196
Connecting to s3-api.us-gio.objectstorage.softlayer.net (s3-api.us-gio.objectstorage.softlayer.net)|67.228.254.196|:443... connecte
HTTP request sent, awaiting response... 200 OK
Length: 3642 (3.6K) [text/csv]
Saving to: 'loan_test.csv'
```

```
loan_test.csv      100%[=====] 3.56K --.-KB/s in 0s
```

```
2023-01-22 11:32:53 (358 MB/s) - 'loan_test.csv' saved [3642/3642]
```

```
test_df = pd.read_csv('loan_test.csv')
test_df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	1	1	PAIDOFF	1000	30	9/8/2016	10/7/2016	50	Bechalor	female
1	5	5	PAIDOFF	300	7	9/9/2016	9/15/2016	35	Master or Above	male
2	21	21	PAIDOFF	1000	30	9/10/2016	10/9/2016	43	High School or Below	female
3	24	24	PAIDOFF	1000	30	9/10/2016	10/9/2016	26	college	male

```
# conversion to datetime object
test_df['due_date'] = pd.to_datetime(test_df['due_date'])
test_df['effective_date'] = pd.to_datetime(test_df['effective_date'])
test_df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender
0	1	1	PAIDOFF	1000	30	2016-09-08	2016-10-07	50	Bechalor	female
1	5	5	PAIDOFF	300	7	2016-09-09	2016-09-15	35	Master or Above	male
2	21	21	PAIDOFF	1000	30	2016-09-10	2016-10-09	43	High School or Below	female

```
# creating weekend column
test_df['dayofweek'] = test_df['effective_date'].dt.dayofweek
test_df['weekend'] = test_df['dayofweek'].apply(lambda x: 1 if (x>3) else 0)
test_df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender	dayofweek
0	1	1	PAIDOFF	1000	30	2016-09-08	2016-10-07	50	Bechalor	female	3
1	5	5	PAIDOFF	300	7	2016-09-09	2016-09-15	35	Master or Above	male	4
2	21	21	PAIDOFF	1000	30	2016-09-10	2016-10-09	43	High School or Below	female	5

```
# replacing male and female string with 0 and 1
test_df['Gender'].replace(to_replace=['male','female'], value=[0,1],inplace=True)
test_df.head()
```

	Unnamed: 0	Unnamed: 0.1	loan_status	Principal	terms	effective_date	due_date	age	education	Gender	dayofweek
0	1	1	PAIDOFF	1000	30	2016-09-08	2016-10-07	50	Bechalor	1	3
1	5	5	PAIDOFF	300	7	2016-09-09	2016-09-15	35	Master or Above	0	4
2	21	21	PAIDOFF	1000	30	2016-09-10	2016-10-09	43	High School or Below	1	5

```
# creaing dummies and dropping one column
X_t = test_df[['Principal','terms','age','Gender','weekend']]
X_t = pd.concat([X_t,pd.get_dummies(test_df['education'])], axis=1)
X_t.drop(['Master or Above'], axis = 1,inplace=True)
X_t.head()
```

	Principal	terms	age	Gender	weekend	Bechalor	High School or Below	college
0	1000	30	50	1	0	1	0	0
1	300	7	35	0	1	0	0	0
2	1000	30	43	1	1	0	1	0
3	1000	30	26	0	1	0	0	1
4	800	15	29	0	1	1	0	0

```
# Feature Scaling
X_t = preprocessing.StandardScaler().fit(X_t).transform(X_t)
X_t[0:5]
```

```
array([[ 0.49362588,  0.92844966,  3.05981865,  1.97714211, -1.30384048,
         2.39791576, -0.79772404, -0.86135677],
       [-3.56269116, -1.70427745,  0.53336288, -0.50578054,  0.76696499,
        -0.41702883, -0.79772404, -0.86135677],
       [ 0.49362588,  0.92844966,  1.88080596,  1.97714211,  0.76696499,
        -0.41702883,  1.25356634, -0.86135677],
       [ 0.49362588,  0.92844966, -0.98251057, -0.50578054,  0.76696499,
        -0.41702883, -0.79772404,  1.16095912],
       [-0.66532184, -0.78854628, -0.47721942, -0.50578054,  0.76696499,
         2.39791576, -0.79772404, -0.86135677]])
```

```
test_df['loan_status'].replace(to_replace=['PAIDOFF', 'COLLECTION'], value=[0,1],inplace=True)
```

```
y_t = test_df['loan_status'].values  
y_t[0:20]
```

```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
```

```
# np arrays to store intermediate result  
Jaccard = np.full(4, np.nan)  
F1_score = np.full(4, np.nan)  
LogLoss = np.full(4, np.nan)  
Algorithm = np.array(4)  
Algorithm = ["KNN", "Decision Tree", "SVM", "LogisticRegression"]
```

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 5:07 PM

