

# Pattern & Anomaly Detection Lab

## Experiment 13

### Random Forest

#### Submitted By:

Dhruv Singhal

500075346

R177219074

AIML B3

#### Submitted To:

Dr. Gopal Phartiyal

Asst. Professor

SOCS

UPES

# CODE:

```
1  from sklearn.datasets import make_classification
2  from sklearn.model_selection import train_test_split
3  from sklearn.ensemble import RandomForestClassifier
4  import seaborn as sns
5  import pandas as pd
6  from sklearn.metrics import confusion_matrix
7  import matplotlib.pyplot as plt
8  ###
9
10 #generate dataset
11 X,y = make_classification(n_samples=1000,n_classes=2)
12 X=pd.DataFrame(X)
13 y=pd.Series(y)
14 print("X values are:",X.head())
15 print("Y values are:",y.head())
16
17 ### Visualization
18
19 plt.hist(X)
20 plt.show()
21
22 ###
23 sns.distplot(y)
24 plt.show()
25
26 ###
27 sns.distplot(X)
28 plt.show()
29
30 ### feature extraction preprocessing, correlation matrix
31 _,graph=plt.subplots(figsize=(15,10))
32 sns.heatmap(X.corr(),annot=True,ax=graph,square=True)
33 plt.show()
34
35
36 ###
37 X_train,X_test,Y_train,Y_test=train_test_split(X,y,test_size=0.15,random_state=42)
38
39 model=RandomForestClassifier()
40 model.fit(X_train,Y_train)
41 print(model.classes_)
42
43 Y_pred=model.predict(X_test)
44 print(Y_pred)
45
46 print("train Accuracy:",model.score(X_train,Y_train))
47 print("test Accuracy:",model.score(X_test,Y_test))
48
49 print(confusion_matrix(Y_test,Y_pred))
50
51
```

# OUTPUT:

```
In [2]: runcell(2, 'B:/3rd year/5th sem/P&AD/random_normal.py')
```

```
X values are:      0      1      2      ...      17      18      19
0  1.370177  0.062732  0.236024  ...  0.109538 -0.201451  2.372252
1  1.454009  1.366763 -0.714880  ...  1.536780 -0.850639  0.993668
2  0.261344 -1.142819 -0.622376  ... -0.733287 -1.585985 -0.511575
3  0.215769 -2.041450 -2.239620  ... -2.029578  0.163686  1.164900
4 -0.654511  0.740525  0.538260  ...  0.251433  1.959272  1.413851
```

```
[5 rows x 20 columns]
```

```
Y values are: 0      1
```

```
1      1
```

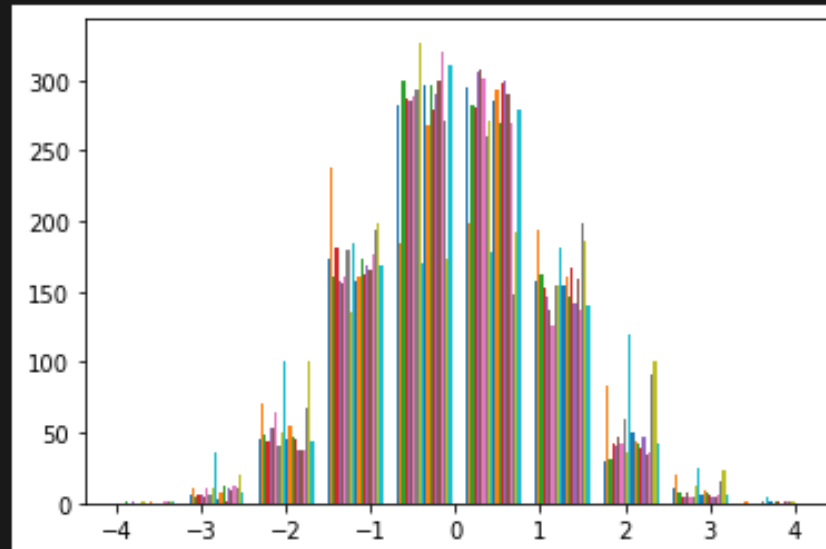
```
2      0
```

```
3      0
```

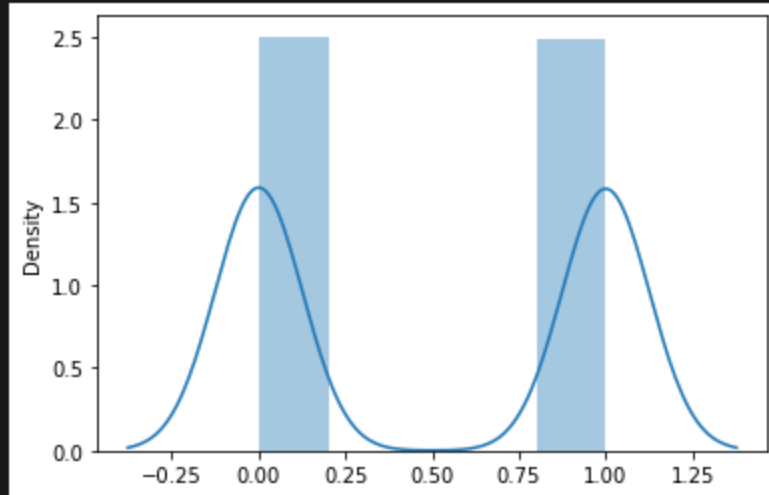
```
4      1
```

```
dtype: int32
```

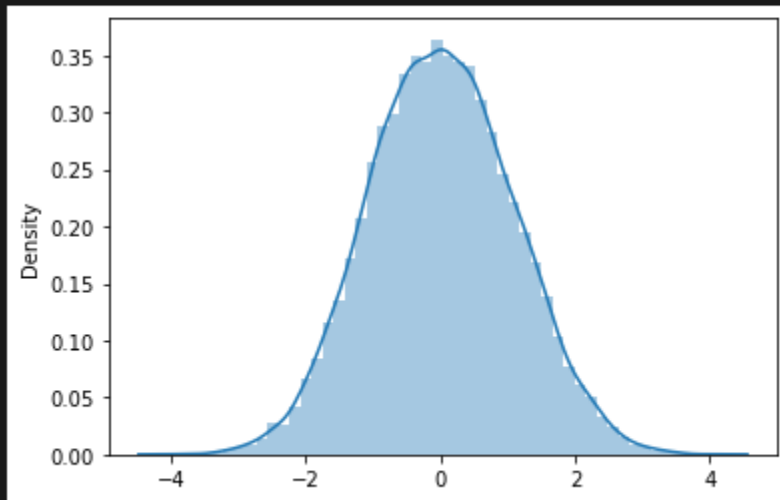
```
In [3]: runcell('Visualization', 'B:/3rd year/5th sem/P&AD/random_normal.py')
```

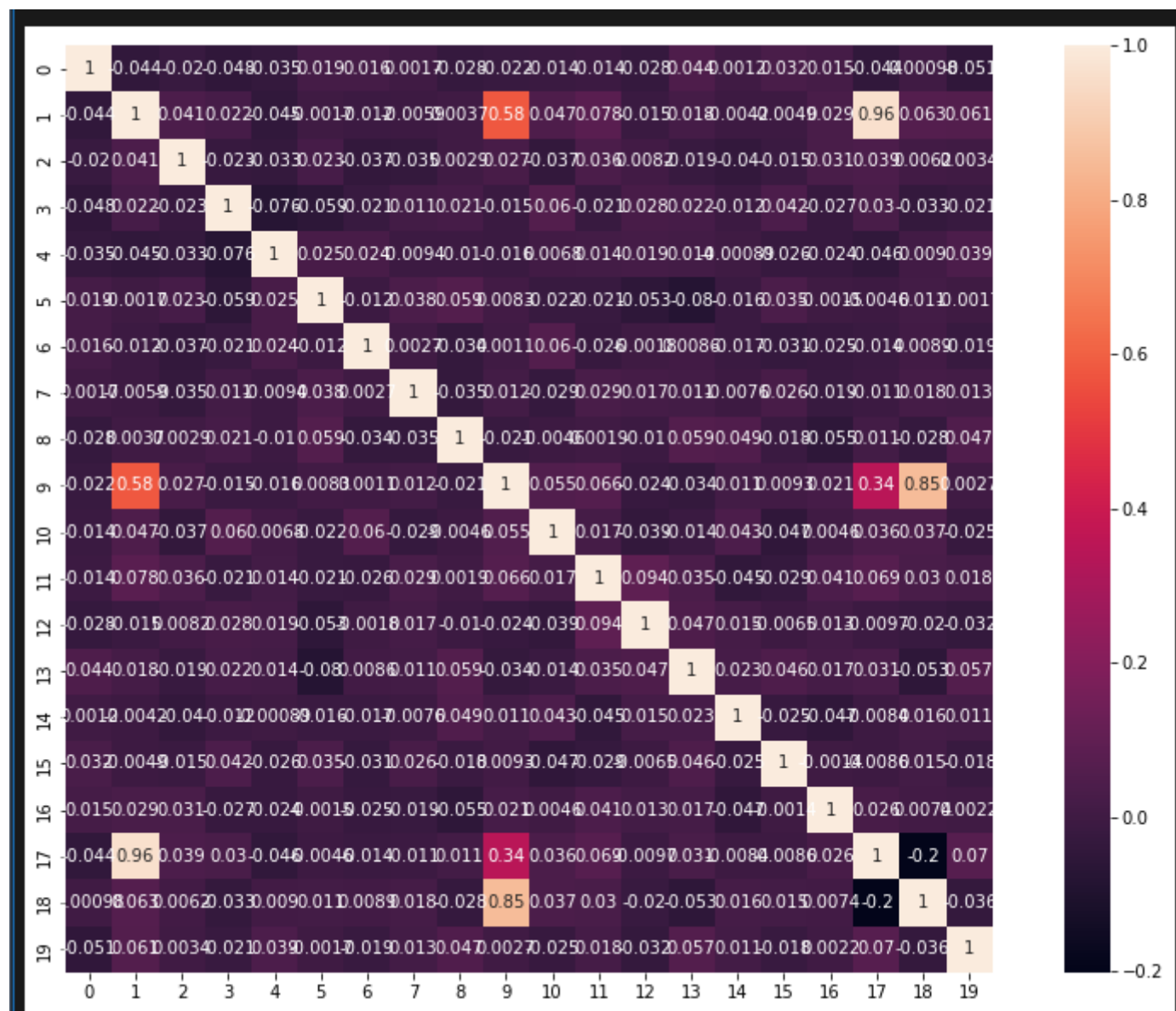


```
In [4]: runcell(4, 'B:/3rd year/5th sem/P&AD/random_normal.py')
C:\Users\Dhruv Singhal\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to
use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level
function for histograms).
  warnings.warn(msg, FutureWarning)
```



```
In [5]: runcell(5, 'B:/3rd year/5th sem/P&AD/random_normal.py')
C:\Users\Dhruv Singhal\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Please adapt your code to
use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level
function for histograms).
  warnings.warn(msg, FutureWarning)
```





```
In [7]: runcell(7, 'B:/3rd year/5th sem/P&AD/random_normal.py')
[0 1]
[0 0 0 1 0 0 0 0 1 1 0 0 0 1 1 0 0 1 0 0 0 0 0 1 0 0 0 0 1 1 0 0 0 0 1 0 0
 1 1 0 0 0 0 1 0 1 1 1 1 0 1 0 1 0 1 1 0 1 1 0 0 1 1 0 1 0 0 0 0 0 0 0 0 0 0
 1 1 1 0 1 0 0 1 1 1 1 1 1 0 0 0 0 1 1 0 1 0 1 0 0 0 1 0 1 0 1 1 1 1 1 0 0
 1 1 0 1 1 1 1 0 0 0 0 1 1 1 0 1 1 0 1 0 1 0 1 0 0 1 1 1 0 1 0 0 1 0 1 1 1
 0 1]
train Accuracy: 1.0
test Accuracy: 0.9133333333333333
[[75  8]
 [ 5 62]]
```

# CODE:

## Multiple Class

```
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import RandomForestClassifier
4 import seaborn as sns
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import warnings
8 warnings.filterwarnings('ignore')
9
10
11 #generate dataset
12 X,y = make_classification(n_samples=1000,n_classes=5,n_features=20,
13 n_informative=10,n_redundant=5,shuffle=True,random_state=42)
14 X=pd.DataFrame(X)
15 y=pd.DataFrame(y)
16 print("X values are:",X.head())
17 print("Y values are:",y.head())
18
19 ### Visualization
20
21 plt.hist(X)
22 plt.show()
23
24 ###
25 sns.distplot(y)
26 plt.show()
27
28 ###
29 sns.distplot(X)
30 plt.show()
31
32 ### feature extractionpreprocessing, correlation matrix
33 _,graph=plt.subplots(figsize=(15,10))
34 sns.heatmap(X.corr(),annot=True,ax=graph,square=True)
35 plt.show()
36 ###
37 X_train,X_test,Y_train,Y_test=train_test_split(X,y,test_size=0.15,random_state=42)
38
39 model=RandomForestClassifier()
40 model.fit(X_train,Y_train)
41 print(model.classes_)
42
43 Y_pred=model.predict(X_test)
44 print(Y_pred)
```



```

45
46 print("train Accuracy:",model.score(X_train,Y_train))
47 print("test Accuracy:",model.score(X_test,Y_test))
48 #%%
49 from sklearn.model_selection import KFold
50 from sklearn.model_selection import cross_val_score
51 clf = RandomForestClassifier(max_depth=5, random_state=0)
52 kf=KFold(n_splits=7)
53 score=cross_val_score(clf, X, y, cv=kf)
54 print("Cross Validation Scores are {}".format(score))
55 print("Average Cross Validation score :{}".format(score.mean()))
56
57
58 #%%
59 from sklearn.model_selection import GridSearchCV
60 tuned_parameters = [{'n_estimators':[10,20,40,100],'criterion':['gini', 'entropy'],
61                       'max_features':['auto', 'sqrt', 'log2'],'bootstrap':[True,False]}]
62 clf=GridSearchCV(RandomForestClassifier(),tuned_parameters,scoring=('accuracy'),verbose=3)
63 clf.fit(X,y)
64 print("Best parameters set found on development set:")
65 print()
66 print(clf.best_params_)
67 print()
68 print("Best Score:",clf.best_score_)
69 z=clf.cv_results_
70

```

# OUTPUT:

```
In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD/exp13_multi.py')
```

```
In [2]: runcell(1, 'B:/3rd year/5th sem/P&AD/exp13_multi.py')
```

X values are:	0	1	2	...	17	18	19
0	8.777957	-1.390182	2.866176	...	-2.372538	0.321516	-3.421634
1	-0.296949	-1.350663	-2.596690	...	0.577353	0.005837	-1.315612
2	3.198753	0.020124	2.434817	...	-1.054403	1.699146	-3.275074
3	4.725085	-0.114722	-0.956705	...	0.660762	-2.022246	-3.007724
4	-3.280166	-0.763541	0.201406	...	-3.444687	-0.181709	-0.402945

```
[5 rows x 20 columns]
```

```
Y values are: 0
```

0 2

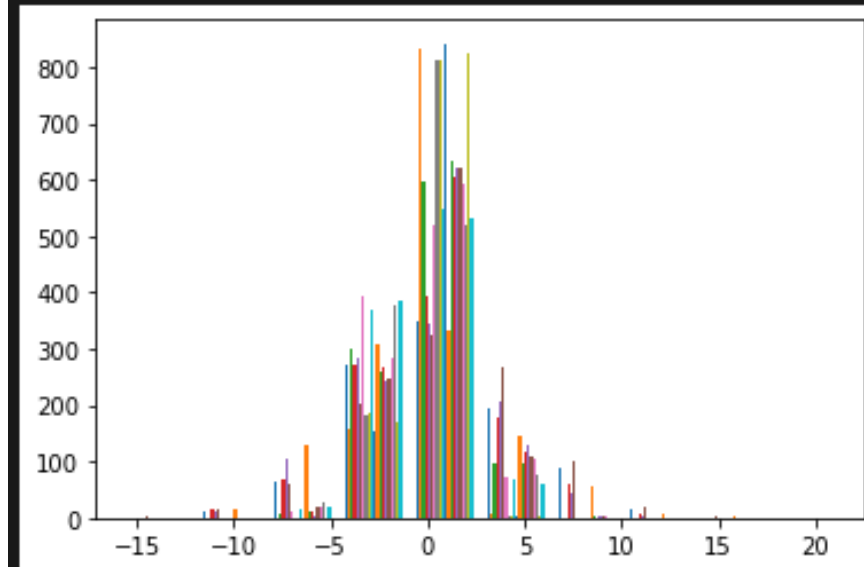
1 3

2 2

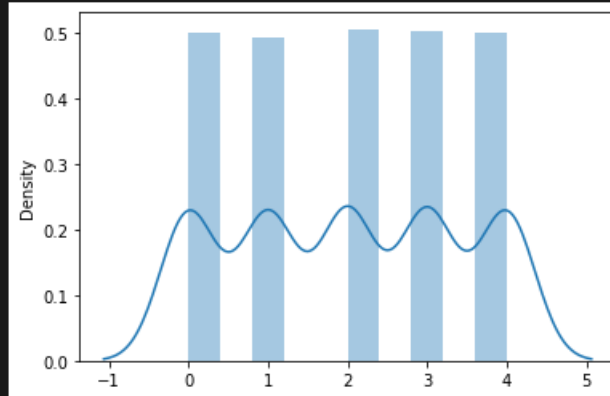
3 0

4 3

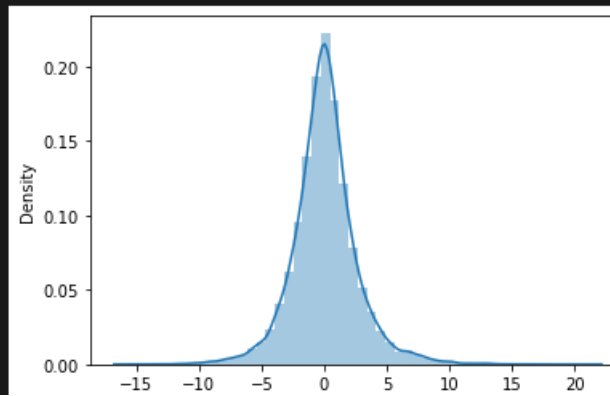
```
In [3]: runcell('Visualization', 'B:/3rd year/5th sem/P&AD/exp13_multi.py')
```



```
In [4]: runcell(3, 'B:/3rd_year/5th sem/P&AD/exp13_multi.py')
```



```
In [5]: runcell(4, 'B:/3rd_year/5th sem/P&AD/exp13_multi.py')
```





```

In [7]: runcell(6, 'B:/3rd year/5th sem/P&AD/exp13_multi.py')
[0 1 2 3 4]
[2 1 2 3 2 2 0 2 0 4 0 1 2 1 1 0 2 4 3 4 2 2 0 1 0 0 3 4 0 0 2 3 3 0 4 1 2
 0 0 1 4 0 1 0 1 1 4 4 4 2 3 3 2 3 1 2 3 3 1 1 1 2 3 0 0 4 0 4 1 3 0 3 1 3
 2 1 0 1 4 0 0 1 1 2 3 2 1 2 0 3 2 3 1 3 4 3 2 4 3 4 1 3 1 4 2 3 0 4 0 0 2
 3 2 2 4 1 1 2 1 3 4 4 2 4 2 1 1 2 2 2 0 2 2 1 3 2 3 0 0 3 0 0 2 4 4 2 4 4
 1 3]
train Accuracy: 1.0
test Accuracy: 0.7333333333333333

In [8]: runcell(7, 'B:/3rd year/5th sem/P&AD/exp13_multi.py')
Cross Validation Scores are [0.6013986  0.62237762 0.62937063 0.64335664 0.6993007  0.6013986
 0.62676056]
Average Cross Validation score :0.6319947657975826

```

Best parameters set found on development set:

```
{'bootstrap': False, 'criterion': 'gini', 'max_features': 'auto', 'n_estimators': 100}
```

Best Score: 0.723