**EXERCISE 2:**

Changes done in "**phishing_sklearn.py**," added f1 score, recall, precision code

```
1.  from warnings import simplefilter
2.
3.  import numpy as np
4.  import pandas as pd
5.  import sklearn
6.  from numpy import genfromtxt
7.  from sklearn import datasets
8.  from sklearn.naive_bayes import GaussianNB
9.  from sklearn.tree import DecisionTreeRegressor
10. from sklearn.ensemble import RandomForestClassifier
11. from sklearn.linear_model import LogisticRegression
12. from sklearn.metrics import (accuracy_score, confusion_matrix, f1_score,
13.                               precision_score, recall_score)
14. from sklearn.model_selection import train_test_split
15. from sklearn.preprocessing import LabelEncoder, StandardScaler
16.
17. simplefilter(action='ignore', category=FutureWarning)
18.
19. ################################################################
20.
21. feature=genfromtxt('phishing.csv',delimiter=',',usecols=(i for i in
    range(0,30)),skip_header=1)
22. target=genfromtxt('phishing.csv',delimiter=',',usecols=(-1),skip_header=1)
23. sc = StandardScaler()
24. sc.fit(feature)
25. target_label = LabelEncoder().fit_transform(target)
26. feature_std = sc.transform(feature)
27. test_size_val=0.33
28. x_train, x_test, y_train, y_test = train_test_split(feature_std, target_label,
    test_size=test_size_val, random_state=1)
29.
30. print("Begin with test_size"+str(test_size_val)+":_____")
31. ##################################################
32. ## print stats
33. precision_scores_list = []
34. accuracy_scores_list = []
35.
36. def print_stats_metrics(y_test, y_pred):
37.     print('Accuracy: %.2f' % accuracy_score(y_test,   y_pred) )
38.     accuracy_scores_list.append(accuracy_score(y_test,   y_pred) )
39.     confmat = confusion_matrix(y_true=y_test, y_pred=y_pred)
40.
41.     print('F1 score: %.2f' % f1_score(y_true=y_test,y_pred=y_pred))
42.     print('Precision: %.2f' % precision_score(y_true=y_test,y_pred=y_pred))
43.     print('Recall: %.2f' % recall_score(y_true=y_test,y_pred=y_pred))
44.
45.     print ("confusion matrix")
46.     print(confmat)
47.     print (pd.crosstab(y_test, y_pred, rownames=['True'], colnames=['Predicted'],
    margins=True))
48.
49. #####################Logistic Regression###########################
50. print("#####################Logistic Regression#############################")
51. clfLog = LogisticRegression()
52. clfLog.fit(x_train,y_train)
53. predictions = clfLog.predict(x_test)
54. print_stats_metrics(y_test, predictions)
55.
56. #####################Random Forest###########################
```

```
57. print("#####################Random Forest##########################")
58. clfRandForest = RandomForestClassifier()
59. clfRandForest.fit(x_train,y_train)
60. predictions = clfRandForest.predict(x_test)
61. print_stats_metrics(y_test, predictions)
62. ####################Decision Tree####################
63. print("####################Decision Tree####################")
64. clfDT = DecisionTreeRegressor()
65. clfDT.fit(x_train,y_train)
66. predictions = clfDT.predict(x_test)
67. print_stats_metrics(y_test, predictions)
68. ####################Naive Bayes####################
69. print("####################Naive Bayes####################")
70. clfNB = GaussianNB()
71. clfNB.fit(x_train,y_train)
72. predictions = clfNB.predict(x_test)
73. print_stats_metrics(y_test, predictions)
```

**At test_size=0.33**

```
1.  ####################Logistic Regression#######################
2.  Accuracy: 0.92
3.  F1 score: 0.93
4.  Precision: 0.92
5.  Recall: 0.94
6.  confusion matrix
7.  [[1476  156]
8.   [ 121 1896]]
9.  Predicted     0     1   All
10. True
11. 0           1476   156  1632
12. 1            121  1896  2017
13. All         1597  2052  3649
14. #####################Random Forest#######################
15. Accuracy: 0.97
16. F1 score: 0.97
17. Precision: 0.96
18. Recall: 0.98
19. confusion matrix
20. [[1557   75]
21.  [  42 1975]]
22. Predicted     0     1   All
23. True
24. 0           1557    75  1632
25. 1             42  1975  2017
26. All         1599  2050  3649
27. #####################Decision Tree####################
28. Accuracy: 0.96
29. F1 score: 0.96
30. Precision: 0.96
31. Recall: 0.96
32. confusion matrix
33. [[1549   83]
34.  [  76 1941]]
35. Predicted    0.0   1.0   All
36. True
37. 0           1549    83  1632
38. 1             76  1941  2017
39. All         1625  2024  3649
40. #####################Naive Bayes####################
41. Accuracy: 0.62
42. F1 score: 0.48
43. Precision: 1.00
```

```
44. Recall: 0.31
45. confusion matrix
46. [[1629    3]
47.  [1383  634]]
48. Predicted    0    1    All
49. True
50. 0           1629    3  1632
51. 1           1383  634  2017
52. All         3012  637  3649
53.
54.
```

**At test_size=0.50**

```
1.  Begin with test_size 0.50:_____
2.  ######################Logistic Regression##########################
3.  Accuracy: 0.92
4.  F1 score: 0.93
5.  Precision: 0.92
6.  Recall: 0.94
7.  confusion matrix
8.  [[2216  248]
9.   [ 176 2888]]
10. Predicted    0    1    All
11. True
12. 0           2216   248  2464
13. 1            176  2888  3064
14. All         2392  3136  5528
15. ######################Random Forest##########################
16. Accuracy: 0.97
17. F1 score: 0.97
18. Precision: 0.96
19. Recall: 0.98
20. confusion matrix
21. [[2331  133]
22.  [  59 3005]]
23. Predicted    0    1    All
24. True
25. 0           2331   133  2464
26. 1             59  3005  3064
27. All         2390  3138  5528
28. ######################Decision Tree####################
29. Accuracy: 0.95
30. F1 score: 0.95
31. Precision: 0.95
32. Recall: 0.96
33. confusion matrix
34. [[2308  156]
35.  [ 125 2939]]
36. Predicted   0.0   1.0    All
37. True
38. 0           2308   156  2464
39. 1            125  2939  3064
40. All         2433  3095  5528
41. ######################Naive Bayes####################
42. Accuracy: 0.61
43. F1 score: 0.47
44. Precision: 0.99
45. Recall: 0.30
46. confusion matrix
47. [[2459    5]
48.  [2134  930]]
49. Predicted    0    1    All
```

```
50. True
51. 0            2459     5  2464
52. 1            2134   930  3064
53. All          4593   935  5528
54.
```

**At test_size=0.20**

```
1.  Begin with test_size0.20:_____
2.  #######################Logistic Regression#########################
3.  Accuracy: 0.92
4.  F1 score: 0.93
5.  Precision: 0.92
6.  Recall: 0.94
7.  confusion matrix
8.  [[ 895  101]
9.   [  68 1147]]
10. Predicted    0     1    All
11. True
12. 0            895   101   996
13. 1             68  1147  1215
14. All          963  1248  2211
15. #######################Random Forest#########################
16. Accuracy: 0.98
17. F1 score: 0.98
18. Precision: 0.97
19. Recall: 0.98
20. confusion matrix
21. [[ 963   33]
22.  [  22 1193]]
23. Predicted    0     1    All
24. True
25. 0            963    33   996
26. 1             22  1193  1215
27. All          985  1226  2211
28. ####################Decision Tree####################
29. Accuracy: 0.96
30. F1 score: 0.96
31. Precision: 0.97
32. Recall: 0.96
33. confusion matrix
34. [[ 955   41]
35.  [  46 1169]]
36. Predicted   0.0   1.0   All
37. True
38. 0            955    41   996
39. 1             46  1169  1215
40. All         1001  1210  2211
41. ######################Naive Bayes####################
42. Accuracy: 0.62
43. F1 score: 0.47
44. Precision: 0.99
45. Recall: 0.31
46. confusion matrix
47. [[993    3]
48.  [843  372]]
49. Predicted    0     1    All
50. True
51. 0            993     3   996
52. 1            843   372  1215
53. All         1836   375  2211
54.
```

**EXERCISE 3:**

**Learning rate:** Itis a tuning parameter in an optimization algorithm that determines the step size at each iteration while moving toward a minimum of a loss function.

**Epoch:** It's a term used in machine learning and indicates the number of passes of the entire training dataset the algorithm has been completed.

**Batch Size:** It's a term used in machine learning and refers to the number of training examples utilized in one iteration.

## Coded F1 scores, precision, recall

```python
1.  import tensorflow.compat.v1 as tf
2.  tf.disable_v2_behavior()
3.  import numpy as np
4.  import pandas as pd
5.  from numpy import genfromtxt
6.  from sklearn import datasets
7.  from sklearn.model_selection import train_test_split
8.  import sklearn
9.  from sklearn.preprocessing import LabelEncoder
10. from sklearn.preprocessing import StandardScaler
11. from sklearn.metrics import accuracy_score
12. from sklearn.metrics import confusion_matrix
13. from sklearn.metrics import precision_score
14. from sklearn.metrics import recall_score, f1_score
15. import pandas as pd
16. import matplotlib.pyplot as plt
17.
18. ################################################################
19.
20. learning_rate = 0.01
21. #n_epochs = 5000
22. n_epochs = 100
23. batch_size = 10000
24.
25. def convertOneHot(data):
26.     y_onehot=[0]*len(data)
27.     for i,j in enumerate(data):
28.         y_onehot[i]=[0]*(data.max()+1)
29.         y_onehot[i][j]=1
30.     return y_onehot
31.
32. ################################################################
33. feature=genfromtxt('phishing.csv',delimiter=',',usecols=(i for i in
    range(0,31)),skip_header=1)
34. target=genfromtxt('phishing.csv',delimiter=',',usecols=(-1),skip_header=1)
35. sc = StandardScaler()
36. sc.fit(feature)
37. target_label = LabelEncoder().fit_transform(target)
38. target_onehot = convertOneHot(target_label)
39. feature_std = sc.transform(feature)
40. x_train, x_test, y_train_onehot, y_test_onehot = train_test_split(feature_std,
    target_onehot, test_size=0.30, random_state=0)
41. A=x_train.shape[1]
42. B=len(y_train_onehot[0])
43. print(len(y_test_onehot[0]))
44. print(B)
```

```python
45. print("Begin:_____")
46. ######################################################################
47.
48. def plot_metric_per_epoch():
49.     x_epochs = []
50.     y_epochs = []
51.     for i, val in enumerate(accuracy_scores_list):
52.         x_epochs.append(i)
53.         y_epochs.append(val)
54.
55.     plt.scatter(x_epochs, y_epochs,s=50,c='lightgreen', marker='s', label='score')
56.     plt.xlabel('epoch')
57.     plt.ylabel('score')
58.     plt.title('Score per epoch')
59.     plt.legend()
60.     plt.grid()
61.     plt.show()
62.
63. ###################################################
64. ## print stats
65. precision_scores_list = []
66. accuracy_scores_list = []
67.
68. def print_stats_metrics(y_test, y_pred):
69.     print('Accuracy: %.2f' % accuracy_score(y_test,   y_pred) )
70.     #Accuracy: 0.84
71.     accuracy_scores_list.append(accuracy_score(y_test,   y_pred) )
72.     confmat = confusion_matrix(y_true=y_test, y_pred=y_pred)
73.
74.     print('F1 score: %.2f' % f1_score(y_true=y_test,y_pred=y_pred))
75.     print('Precision: %.2f' % precision_score(y_true=y_test,y_pred=y_pred))
76.     print('Recall: %.2f' % recall_score(y_true=y_test,y_pred=y_pred))
77.
78.     print ("confusion matrix")
79.     print(confmat)
80.     print (pd.crosstab(y_test, y_pred, rownames=['True'], colnames=['Predicted'],
    margins=True))
81.
82. ###########################################################
83. def layer(input, weight_shape, bias_shape):
84.     weight_stddev = (2.0/weight_shape[0])**0.5
85.     w_init = tf.random_normal_initializer(stddev=weight_stddev)
86.     bias_init = tf.constant_initializer(value=0)
87.     W = tf.get_variable("W", weight_shape, initializer=w_init)
88.     b = tf.get_variable("b", bias_shape, initializer=bias_init)
89.     return tf.nn.relu(tf.matmul(input, W) + b)
90. ###########################################################
91. def inference_deep_layers(x_tf, A, B):
92.     with tf.variable_scope("hidden_1"):
93.         hidden_1 = layer(x_tf, [A, 15],[15])
94.     with tf.variable_scope("hidden_2"):
95.         hidden_2 = layer(hidden_1, [15, 5],[5])
96.     with tf.variable_scope("output"):
97.         output = layer(hidden_2, [5, B], [B])
98.     return output
99. ###########################################################
100.  def loss_deep(output, y_tf):
101.      xentropy = tf.nn.softmax_cross_entropy_with_logits(logits=output, labels=y_tf)
102.      loss = tf.reduce_mean(xentropy)
103.      return loss
104.  ###########################################################
105.
106.  def training(cost):
107.      optimizer = tf.train.GradientDescentOptimizer(learning_rate)
108.      train_op = optimizer.minimize(cost)
```

```
109.        return train_op
110.
111.    #############################################################
112.    def evaluate(output, y_tf):
113.        correct_prediction = tf.equal(tf.argmax(output,1), tf.argmax(y_tf,1))
114.        accuracy = tf.reduce_mean(tf.cast(correct_prediction, "float"))
115.        return accuracy
116.    #############################################################
117.
118.    x_tf = tf.placeholder("float",[None,A])
119.    y_tf = tf.placeholder("float",[None,B])
120.    #############################################################
121.    output = inference_deep_layers(x_tf,A,B)
122.    cost = loss_deep(output,y_tf)
123.    train_op=training(cost)
124.    eval_op=evaluate(output,y_tf)
125.    #############################################################
126.    init = tf.global_variables_initializer()
127.    sess = tf.Session()
128.    sess.run(init)
129.    #############################################################
130.    y_p_metrics = tf.argmax(output,1)
131.    #############################################################
132.    num_samples_train_set=x_train.shape[0]
133.    num_batches = int(num_samples_train_set/batch_size)
134.
135.    #############################################################
136.
137.    for i in range(n_epochs):
138.        print("epoch %s out of %s"%(i,n_epochs))
139.        for batch_n in range(num_batches):
140.            sta = batch_n*batch_size
141.            end = sta+batch_size
142.            sess.run(train_op,feed_dict={x_tf:x_train[sta:end],y_tf:y_train_onehot[sta:end]})
143.        print ("------------------------------------------------------------------------
    -")
144.        print ("Accuracy score")
145.        #result = sess.run(eval_op,feed_dict={x_tf:x_test,y_tf:y_test_onehot})
146.        result, y_result_metrics = sess.run([eval_op, y_p_metrics], feed_dict={x_tf: x_test,
    y_tf: y_test_onehot})
147.        print("Run {},{}".format(i,result))
148.        y_true = np.argmax(y_test_onehot,1)
149.        print_stats_metrics(y_true, y_result_metrics)
150.        if i==n_epochs-1:
151.            plot_metric_per_epoch()
```

## At different epoch, batch size and learning rate

### 1) At

```
1. learning_rate = 0.01
2. n_epochs = 100
3. batch_size = 10000
```
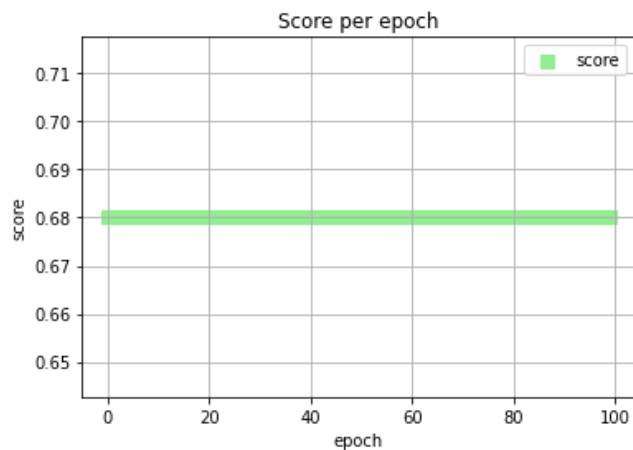
**Output:**

```
1.  Accuracy score
2.  Run 99,0.4727163016796112
3.  Accuracy: 0.47
4.  F1 score: 0.09
5.  Precision: 0.85
6.  Recall: 0.05
7.  confusion matrix
8.  [[1483   15]
9.   [1734   85]]
10. Predicted     0    1   All
11. True
12. 0            1483   15  1498
13. 1            1734   85  1819
14. All          3217  100  3317
15.
16.
```



Score per epoch

2) At

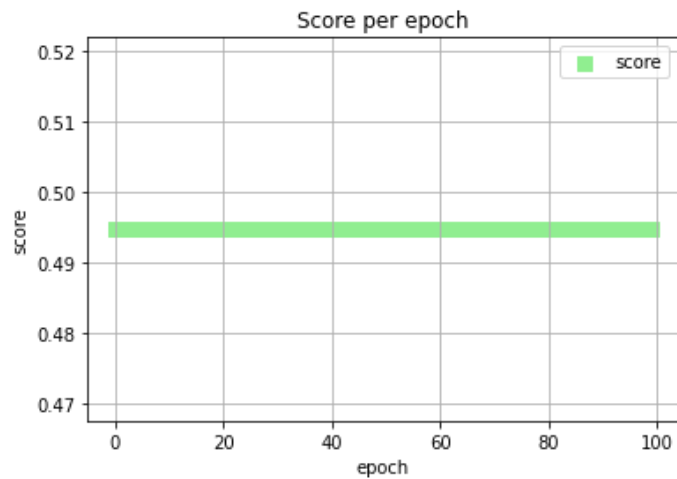```
1.  learning_rate = 0.006
2.  n_epochs = 100
3.  batch_size = 10000
```

**Output:**

```
1.  Accuracy: 0.41
2.  F1 score: 0.20
3.  Precision: 0.39
4.  Recall: 0.13
5.  confusion matrix
6.  [[1127  371]
7.   [1579  240]]
8.  Predicted     0    1   All
9.  True
10. 0            1127  371  1498
11. 1            1579  240  1819
12. All          2706  611  3317
13.
```
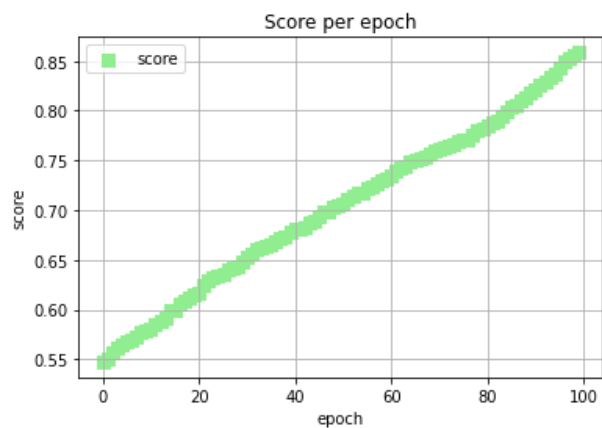
Score per epoch

**3) At**

```
4.  learning_rate = 0.006
5.  n_epochs = 100
6.  batch_size = 1000
```

**Output:**

```
1.  Accuracy: 0.95
2.  F1 score: 0.95
3.  Precision: 0.94
4.  Recall: 0.96
5.  confusion matrix
6.  [[1390  108]
7.   [  65 1754]]
8.  Predicted      0      1    All
9.  True
10. 0            1390    108  1498
11. 1              65   1754  1819
12. All          1455   1862  3317
13.
```



Score per epoch

**4) At**

```
7.   learning_rate = 0.006
8.   n_epochs = 100
9.   batch_size = 128
```
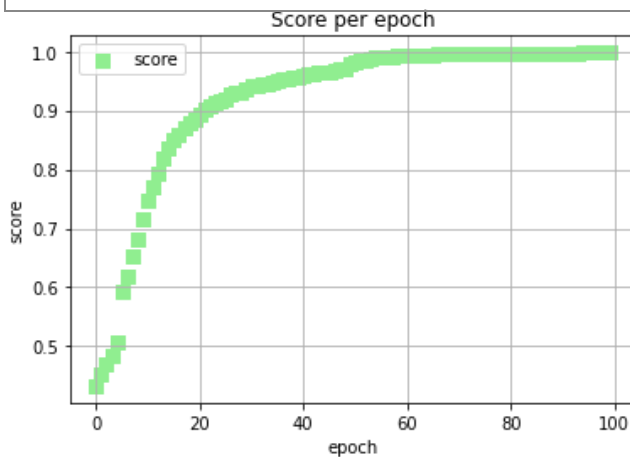
**Output:**

```
1.   Accuracy: 0.99
2.   F1 score: 0.99
3.   Precision: 0.99
4.   Recall: 1.00
5.   confusion matrix
6.   [[1474   24]
7.    [   0 1819]]
8.   Predicted     0      1    All
9.   True
10. 0            1474     24  1498
11. 1               0   1819  1819
12. All          1474   1843  3317
13.
```


Score per epoch

Results:

At first we saw that how our dataset looks like, Then we installed weka and opened the CSV data set. Using the Explorer window and. In the classifier tab, we set the classification as different sets. Logistic regression, Naïve Bayes, Random Forest, etc.

Will logged the data by each of them, The output contains the time taken by the model. The confusion matrix the F1 score, precision, accuracy and recall, for the data set over the trained model.

Time taken by model:

Random Forest: 2.21 s

LMT: 6.89 s

Logistic: 0.3

SGD: 0.49

Naïve bayes: 0.06

Mean Abs Error:

Random Forest: 0.049

LMT: 0.034
Logistic:  0.1078
SGD: 0.0705
Naïve bayes: 0.1197

Precision:
Random Forest: 0.970
LMT: 0.967
Logistic:  0.926
SGD: 0.930
Naïve bayes: 0.909

On getting such results we can finally discuss which phising detector performed better in comparison to others, we can see performance on basis of various parameters that Random Forest performed better in comparison to others as it have high precision score and lower Mean Abs Error, but time taken is much higher compared to rest algorithms

In sklearn we implemented the classifier model on the basis of different algorithms like Logistic Regression, Decision Tree, Random Forest, and Naïve Bayes

On the basis of various parameters and accuracy we can see that which model performed better
Accuracy
Logistic Regression:0.92
Random Forest:0.97
Decision Tree:0.96
Naïve Bayes:0.62
F1 score
Logistic Regression:0.93
Random Forest:0.97
Decision Tree:0.96
Naïve Bayes:0.48
Here also we can see the Random Forest worked better in comparison to other

Now, in terms of NN layer, we saw the highest accuracy came out to be 0.99 on slower learning rate and smaller batch size which was even better than earlier attempts.
So we can say that Neural Networks worked better in comparison to others.
We can also say that Neural Networks have the ability to learn and model non-linear and complex relationships , which is really important because in real-life, many of the relationships between inputs and outputs are non-linear as well as complex, so it performed better.