

1 ML Metrics

1.1 ROC

1. We want to differentiate Eve from Bob. In monitoring Eve's accesses to email on a daily basis we find daily accesses of 20, 25, 16, 42 and 22, and then monitor Bob's accesses as: 50, 41, 60, 54 and 39. With the following we aim to detect Bob from Eve, and plot the ROC Curve. Use the following code to determine the ROC curve and the AUC value:

```
1 https://repl.it/@billbuchanan/class01

21 # https://asecuritysite.com/bigdata/roc
22 from sklearn import metrics
23 import matplotlib.pyplot as plt
24
25 def show_roc(FPR, TPR, AUC):
26     plt.plot(FPR, TPR, color='blue', label='ROC')
27     plt.plot([0, 1], [0, 1], color='darkblue', linestyle='--')
28     plt.xlabel('FPR')
29     plt.ylabel('TPR')
30     plt.title('Receiver Operating Characteristic (ROC) Curve')
31     plt.legend(["AUC=%.3f" % AUC])
32
33     plt.show()
34
35 y = ['Eve', 'Eve', 'Eve', 'Eve', 'Eve', 'Bob', 'Bob', 'Bob', 'Bob', 'Bob']
36 scores = [20,25,16,42,22,50,41,60,54,39]
37
38 positive_label = 'Bob'
39
40 fpr, tpr, thresholds = metrics.roc_curve(y, scores, pos_label=
41     positive_label)
42
43 auc=metrics.auc(fpr, tpr)
44
45 print ("FPR:",fpr)
46 print ("TPR:",tpr)
47
48 print ("Thresholds:",thresholds)
49
50 show_roc(fpr, tpr, auc)
```

What is the AUC:

What are the thresholds used?

For each threshold, what is the FPR and what is the TPR:

If we set a threshold of 42, what is the FPR and what is the TPR?

Bob's daily accesses for email are now monitored for 50, 55, 43, 90, 110 and 66, and Eve has accesses of 14, 32, 19, 46, 21, 48 and 50. Use the program to determine the new ROC curve:

What is the AUC:

What are the thresholds used?

For each threshold, what is the FPR and what is the TPR:

If we set a threshold of 42, what is the FPR and what is the TPR?

2. Now, we can add Alice, and who has accesses of 13, 23, 32, 40, 11, and 14, and determine the following:

What is the AUC:

What are the thresholds used?

For each threshold, what is the FPR and what is the TPR:

If we set a threshold of 42, what is the FPR and what is the TPR?

3. Now change the positive label to Alice, and determine the following:

What is the AUC:

What are the thresholds used?

For each threshold, what is the FPR and what is the TPR:

If we set a threshold of 42, what is the FPR and what is the TPR?

4. In the following example we will load a dataset for a machine learning model to differentiate hand written digits. Run the following code and determine the confusion matrix:

```
1 https://repl.it/@billbuchanan/class02

2 # https://asecuritysite.com/bigdata/sk01
3 import sys
4 import matplotlib.pyplot as plt
5
6 ga=0.011
7
8 from sklearn import datasets, svm, metrics
9
10 digits = datasets.load_digits()
11
12 images_and_labels = list(zip(digits.images, digits.target))
13 for index, (image, label) in enumerate(images_and_labels[:10]):
14     plt.subplot(2, 10, (index + 1))
15     plt.axis('off')
16     plt.imshow(image, cmap=plt.cm.gray_r, interpolation='nearest')
17     plt.title('Tr: %i' % label)
18
19 # To apply a classifier on this data, we need to flatten the image, to
20 # turn the data in a (samples, feature) matrix:
21 n_samples = len(digits.images)
22 data = digits.images.reshape((n_samples, -1))
23
24 # Create a classifier: a support vector classifier
25 classifier = svm.SVC(gamma=ga)
26
27 # We learn the digits on the first half of the digits
28 classifier.fit(data[:n_samples // 2], digits.target[:n_samples // 2])
29
30 # Now predict the value of the digit on the second half:
31 expected = digits.target[n_samples // 2:]
32 predicted = classifier.predict(data[n_samples // 2:])
33
34 print("Classification report for classifier %s:\n%s\n"
35       % (classifier, metrics.classification_report(expected, predicted)))
36 print("Confusion matrix:\n%s" % metrics.confusion_matrix(expected,
37                  predicted))
38
39 images_and_predictions = list(zip(digits.images[n_samples // 2:],
40                                  predicted))
41 for index, (image, prediction) in enumerate(images_and_predictions[:4]):
```

```

39 plt.subplot(2, 4, index + 5)
40 plt.axis('off')
41 plt.imshow(image, cmap=plt.cm.gray_r, interpolation='nearest')
42 plt.title('Prediction: %i' % prediction)
43
44 plt.show()

```

What are the TP for the character '0'?

What are the FP for the character '0'?

What are the FN for the character '0'?

We are using SVM (Support Vector Machine), and which uses a gamma factor. Vary the gamma value with 0.1, 0.2, 0.3 and so on, up to 1.0, and observe how the confusion matrix changes:

5. The following code generates a data set which has two clusters, and then marks each of the dataset elements for their cluster source. Run the program several times and observe the creation of the clusters:

```

1 https://repl.it/@billbuchanan/class03

```

```

1 # This code creates a data set with two clusters (defined by the two
   features. The output is then data_vals[:, 0] and data_vals[:, 1] and
   these are marked by class_lab
2
3 from sklearn.datasets import make_classification
4 import matplotlib.pyplot as plt
5
6 data_vals, class_label = make_classification(n_samples=100, n_features=2,
      n_redundant=0, n_informative=1, n_clusters_per_class=1)
7
8 plt.scatter(data_vals[:, 0], data_vals[:, 1], marker='o', c=class_label, s
      =25, edgecolor='k')
9
10 plt.savefig('test.png')
11
12 plt.show()

```

Modify the code so that it now generates 250 points.

6. We will now use this method of cluster generation, and then split the data into 70% training data, and 30% test data, in to train a RandomForestClassifier model to predict the results:

```
1 https://repl.it/@billbuchanan/class04

1 # Create ROC
2 from sklearn.datasets import make_classification
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import roc_auc_score
6 from sklearn.metrics import roc_curve
7 import matplotlib.pyplot as plt
8
9 def show_roc(FPR, TPR, AUC):
10     plt.plot(FPR, TPR, color='blue', label='ROC')
11     plt.plot([0, 1], [0, 1], color='darkblue', linestyle='--')
12     plt.xlabel('FPR')
13     plt.ylabel('TPR')
14     plt.title('Receiver Operating Characteristic (ROC) Curve')
15     plt.legend(["AUC=%.3f" % AUC])
16     plt.show()
17
18 data_vals, class_label =make_classification(n_samples=10,n_features=2,
19     n_redundant=0, n_informative=1,n_clusters_per_class=1)
20 X_train, X_test, y_train, y_test = train_test_split(data_vals,
21     class_label, test_size=0.3, random_state=1)
22 model = RandomForestClassifier()
23
24 model.fit(X_train, y_train)
25 print ("Model score: ",model.score(X_test, y_test))
26
27 probs = model.predict_proba(X_test)
28
29 # probabilities of getting a 1
30 probs = probs[:, 1]
31
32 auc = roc_auc_score(y_test, probs)
33
34 FPR, TPR, thresholds = roc_curve(y_test, probs)
35
36 print ("Thresholds: ",thresholds)
37 print ("FPR: ",FPR)
38 print ("TPR: ",TPR)
39
40 show_roc(FPR, TPR, auc)
41
```

```
42 plt.scatter(data_vals[:, 0], data_vals[:, 1], marker='o', c=class_labels,
43             =25, edgecolor='k')
44 plt.savefig('test.png')
```

For 10 points, what is the AUC?

For 100 points, what is the AUC?

For 250 points, what is the AUC?

For 1000 points, what is the AUC?

For the different number of points, how does the shape of ROC Curve change?

7. There are a few ensemble methods for machine learning in skLearn, including BaggingClassifier, ExtraTreesClassifier, AdaBoostClassifier, and GradientBoostingRegressor. Modify the code given in Q.6 to support each of the different models:

```
1 from sklearn.ensemble import AdaBoostClassifier
2 ...
3 model = AdaBoostClassifier()
```

For each of the methods, what is the AUC, and which method is the best performing?

ExtraTreesClassifier AUC:

AdaBoostClassifier AUC:

GradientBoostingRegressor AUC:

BaggingClassifier AUC:

8. You have been asked to identify if there is a linkage between gun ownership and population density in US states, and whether there is a link to the number of murders per 100K of the population. An outline of the code is given here:

```
1 https://repl.it/@billbuchanan/class08
```

```

1 # Numeric Prediction
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestRegressor
5 from sklearn.metrics import r2_score, mean_squared_error
6
7 # Features
8 x1= "Gun ownership"
9 x2= "Population density"
10 # Prediction
11 x3 = "Murders per 100K"
12
13 fdata="guns.csv"
14
15 print ("Training data:\t\t",x1,",",x2)
16 print ("Training against:\t",x3)
17 print ("Data set:\t\t",fdata)
18
19 ver=pd.read_csv(fdata)
20
21 dataset=ver[[x1,x2]]
22 train=ver[x3]
23 print (dataset)
24
25 x_train, x_test, y_train, y_test= train_test_split(dataset,train,
26                                                     test_size=0.3, random_state=1)
27
28 model= RandomForestRegressor()
29 model.fit(x_train,y_train)
30
31 y_predictions =model.predict(x_test)
32
33 accuracy = r2_score(y_test, y_predictions)
34 mse = mean_squared_error(y_test, y_predictions)
35
36
37 print ("R^2=",accuracy)
38 print ("MSE=",mse)
39
40 # Correlation
41 cor=ver.corr()
42 print (cor[x3])

```

By examining the R^2 value, is the machine learning implementation a good model?

Can you create a better model for predicting Murders per 100K of the population, and with only two features?