

EE2211 Tutorial 10

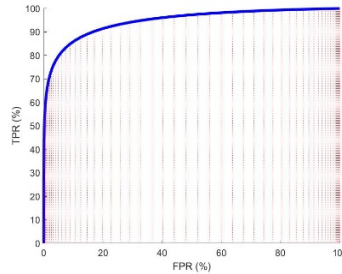
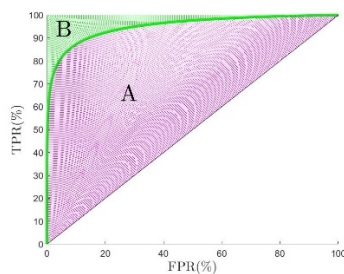
Question 1:

We have two classifiers showing the same accuracy with the same cross-validation. The more complex model (such as a 9th-order polynomial model) is preferred over the simpler one (such as a 2nd-order polynomial model).

- a) True
- b) False

Answer: b).

Question 2: According to the plots below, the Gini Coefficient is equal to Two times the Area Under the ROC minus One.



- a) True
- b) False

Answer: a).

Reason: Since the area $(A+B) = \frac{1}{2}$ (half the square box of area 1), Gini-coefficient = $A/(A+B) = 2A$.

$AUC = A + 1/2 \Rightarrow A = AUC - 0.5$. Substitute A into the Gini above: Gini-coefficient = $2(AUC - 0.5)$.

Question 3:

Suppose the binary classification problem, which you are dealing with, has highly imbalanced classes. The majority class has 99 hundred samples and the minority class has 1 hundred samples. Which of the following metric(s) would you choose for assessing the classification performance? (Select all relevant metric(s) to get full credit)

- a) Classification Accuracy
- b) Cost sensitive accuracy
- c) Precision and recall
- d) None of these

Answer: (b, c)

Question 4:

Given below is a scenario for Training error rate Tr , and Validation error rate Va for a machine learning algorithm. You want to choose a hyperparameter (P) based on Tr and Va .

P	Tr	Va
---	----	----

10	0.10	0.25
9	0.30	0.35
8	0.22	0.15
7	0.15	0.25
6	0.18	0.15

Which value of P will you choose based on the above table?

- a) 10
- b) 9
- c) 8
- d) 7
- e) 6

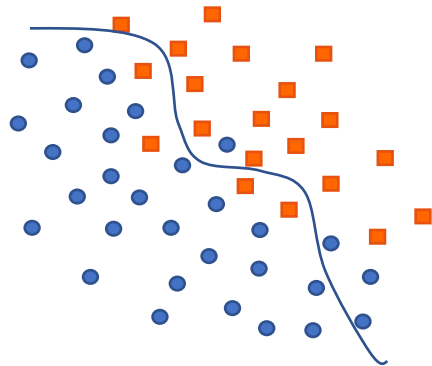
Answer: e).

(Binary and Multicategory Confusion Matrices)

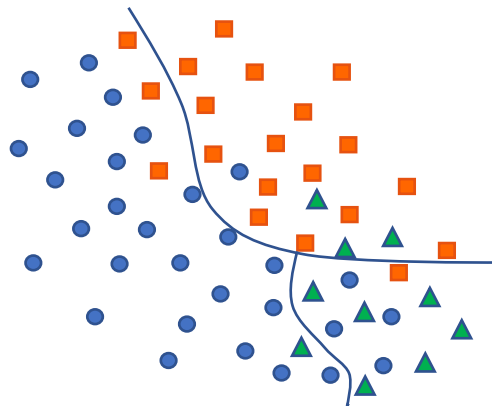
Question 5:

Tabulate the confusion matrices for the following classification problems.

(a) Binary problem (the class-1 and class-2 data points are respectively indicated by squares and circles)



(b) Three-category problem (the class-1, class-2 and class-3 data points are respectively indicated by squares, circles and triangles).



Answer:

(a)

	P_1	P_2
P_1	16	4
P_2	4	26

(b)

	P_1	P_2	P_3
P_1	16	3	1
P_2	1	25	4
P_3	3	1	6

(5-fold Cross-validation)

Question 6:

Get the data set “from sklearn.datasets import load_iris”. Perform a 5-fold Cross-validation to observe the best polynomial order (among orders 1 to 10 and without regularization) for validation prediction. Note that, you will have to partition the whole dataset for training/validation/test parts, where the size of validation set is the same as that of test. Provide a plot of the average 5-fold training and validation error rates over the polynomial orders. The randomly partitioned data sets of the 5-fold shall be maintained for reuse in evaluation of future algorithms.

Answer:

```
##--- load data from scikit ---##
import numpy as np
import pandas as pd
print("pandas version: {}".format(pd.__version__))
import sklearn
print("scikit-learn version: {}".format(sklearn.__version__))
from sklearn.datasets import load_iris
iris_dataset = load_iris()
X = np.array(iris_dataset['data'])
y = np.array(iris_dataset['target'])
## one-hot encoding
Y = list()
for i in y:
    letter = [0, 0, 0]
    letter[i] = 1
    Y.append(letter)
Y = np.array(Y)
test_idx = np.random.RandomState(seed=2).permutation(Y.shape[0])
X_test = X[test_idx[:25]]
Y_test = Y[test_idx[:25]]
X = X[test_idx[25:]]
Y = Y[test_idx[25:]]

from sklearn.preprocessing import PolynomialFeatures
error_rate_train_array = []
error_rate_val_array = []
##--- Loop for Polynomial orders 1 to 10 ---##
for order in range(1,11):
```

```

error_rate_train_array_fold = []
error_rate_val_array_fold = []
# Random permutation of data
Idx = np.random.RandomState(seed=8).permutation(Y.shape[0])
# Loop 5 times for 5-fold
for k in range(0,5):
    ##--- Prepare training, validation, and test data for the 5-fold ---#
    # Prepare indexing for each fold
    X_val = X[Idx[k*25:(k+1)*25]]
    Y_val = Y[Idx[k*25:(k+1)*25]]
    Idxtrn = np.setdiff1d(Idx, Idx[k*25:(k+1)*25])
    X_train = X[Idxtrn]
    Y_train = Y[Idxtrn]
    ##--- Polynomial Classification ---##
    poly = PolynomialFeatures(order)
    P = poly.fit_transform(X_train)
    Pval = poly.fit_transform(X_val)
    if P.shape[0] > P.shape[1]: # over-/under-determined cases
        reg_L = 0.00*np.identity(P.shape[1])
        inv_PTP = np.linalg.inv(P.transpose().dot(P)+reg_L)
        pinv_L = inv_PTP.dot(P.transpose())
        wp = pinv_L.dot(Y_train)
    else:
        reg_R = 0.00*np.identity(P.shape[0])
        inv_PPT = np.linalg.inv(P.dot(P.transpose())+reg_R)
        pinv_R = P.transpose().dot(inv_PPT)
        wp = pinv_R.dot(Y_train)
    ##--- trained output ---##
    y_est_p = P.dot(wp);
    y_cls_p = [[1 if y == max(x) else 0 for y in x] for x in y_est_p ]
    m1tr = np.matrix(Y_train)
    m2tr = np.matrix(y_cls_p)
    # training classification error count and rate computation
    difference = np.abs(m1tr - m2tr)
    error_train = np.where(difference.any(axis=1))[0]
    error_rate_train = len(error_train)/len(difference)
    error_rate_train_array_fold += [error_rate_train]
    ##--- validation output ---##
    yval_est_p = Pval.dot(wp);
    yval_cls_p = [[1 if y == max(x) else 0 for y in x] for x in yval_est_p ]
    m1 = np.matrix(Y_val)
    m2 = np.matrix(yval_cls_p)
    # validation classification error count and rate computation
    difference = np.abs(m1 - m2)
    error_val = np.where(difference.any(axis=1))[0]
    error_rate_val = len(error_val)/len(difference)
    error_rate_val_array_fold += [error_rate_val]
    # store results for each polynomial order
    error_rate_train_array += [np.mean(error_rate_train_array_fold)]
    error_rate_val_array += [np.mean(error_rate_val_array_fold)]
##--- plotting ---##
import matplotlib.pyplot as plt
order=[x for x in range(1,11)]
plt.plot(order, error_rate_train_array, color='blue', marker='o', linewidth=3,
label='Training')
plt.plot(order, error_rate_val_array, color='orange', marker='x', linewidth=3,
label='Validation')
plt.xlabel('Order')
plt.ylabel('Error Rates')
plt.title('Training and Validation Error Rates')
plt.legend()
plt.show()

```



(Plot the ROC and Compute the AUC for Binary Classification)

Question 7:

Download the spambase data set from the UCI Machine Learning repository

<https://archive.ics.uci.edu/ml/machine-learning-databases/spambase/> and use the following function to pack the data:

```
def load_data(Train=False):
    import csv
    data = []
    ## Read the training data
    f = open('spambase.data')
    reader = csv.reader(f)
    next(reader, None)
    for row in reader:
        data.append(row)
    f.close()
    ## x[:-1]: omit the last element of each x row
    X = np.array([x[:-1] for x in data]).astype(np.float)
    ## x[-1]: the first element from the right instead of from the left
    y = np.array([x[-1] for x in data]).astype(np.float)
    del data # free up the memory
    if Train:
        # returns X_train, X_test, y_train, y_test
        return train_test_split(X, y, test_size=0.2, random_state=8)
    else:
        return X, y
```

Randomly split the dataset into two parts, 80% for training and 20% for testing. Compute the test Classification Error Rate and the AUC based on the optimal linear regression model without regularization.

In other words, use the training set to train a linear model and use the test set to check the classification performance in terms of Classification Error Rate, ROC and AUC.

Hint: to plot the ROC curve, the population of output predictions (say, values stored in vector `y_predict`) needs to be separated according to the two known output classes (0 or 1 values in the target vector `y_test`). Let `y_predict_for_PosSamples` (when `y_test==1`) and `y_predict_for_NegSamples` (when `y_test==0`) denote these two populations of prediction. Then compute the TPR (=1-FNR) and the FPR at various threshold/operating points in order to plot the ROC curve. To obtain the highest possible resolution for the ROC plot, you can set the decision threshold according to each element of the lower population of the two predicted classes of data (`y_predict_for_PosSamples`, `y_predict_for_NegSamples`).

Answer

```
import numpy as np
from numpy.random import RandomState
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt

## Get training and test sets
X_train, X_test, y_train, y_test = load_data(Train=True)
## Linear regression
inv_XTX = np.linalg.inv(X_train.transpose().dot(X_train))
pinv = inv_XTX.dot(X_train.transpose())
W = pinv.dot(y_train)
## Prediction
y_predict = X_test.dot(W)
## Calculate classification error rate
yp_cls = [1 if yout >= 0.5 else 0 for yout in y_predict]
difference = np.abs(y_test - yp_cls)
test_error_count = (difference == 1).sum()
test_error_rate = test_error_count/len(y_test)
print(test_error_rate)
##--- Compute FPR and FNR at different thresholds ---##
## separate the two classes of predicted data based on the ground truth y_test
pos_idx = np.where(y_test == 1) # identify the indexing of positive-class in the test set
neg_idx = np.where(y_test == 0) # identify the indexing of negative-class in the test set
y_predict_for_PosSamples = y_predict[pos_idx] # prediction of the positive-class data
y_predict_for_NegSamples = y_predict[neg_idx] # prediction of the negative-class data
## use the shorter among the two arrays as threshold
if ( len(y_predict_for_PosSamples) <= len(y_predict_for_NegSamples) ):
    sorted = np.sort(y_predict_for_PosSamples) # sort in ascending order to be used as threshold
else:
    sorted = np.sort(y_predict_for_NegSamples) # sort in ascending order to be used as threshold
FNR = []
FPR = []
TPR = []
## Compute FNR, FPR, and TPR for each threshold
for k in range(len(sorted)):
    yp_cls_pos = np.abs([1 if yout >= sorted[k] else 0 for yout in y_predict_for_PosSamples])
    yp_cls_neg = np.abs([1 if yout >= sorted[k] else 0 for yout in y_predict_for_NegSamples])
    FNR += [(yp_cls_pos == 0).sum()/len(y_predict_for_PosSamples)]
    FPR += [(yp_cls_neg == 1).sum()/len(y_predict_for_NegSamples)]
    TPR += [1-(yp_cls_pos == 0).sum()/len(y_predict_for_PosSamples)]
##--- Plot ROC and DET curves ---##
plt.plot(FPR, FNR, '-', label = 'DET')
plt.plot(FPR, TPR, '-', label = 'ROC')
plt.xlabel('FPR')
plt.ylabel('FNR/TPR')
plt.legend(fontsize=15)
##--- Compute AUC ---##
ypos_array = [[1 if y_predict_for_PosSamples[j] >= y_predict_for_NegSamples[k] else 0 for j in
range(len(y_predict_for_PosSamples))] for k in range(len(y_predict_for_NegSamples))]
AUC = np.sum(ypos_array)/(len(y_predict_for_PosSamples)*len(y_predict_for_NegSamples))
print(AUC)
```

Test Error Rate and AUC:

0.15108695652173912

0.9432018448840879

