

```

# Assessed exercises 7
# Look at cuts and creating ROC curves

from pandas import Series, DataFrame
import pandas as pd
import numpy as np
import numpy.random as npr
from matplotlib import pyplot as plt

# This week we will use the ebola_test dataset to test the code. It contains the
# results for a new medical test for detecting Ebola that is much faster, though
# less accurate, than that currently available. The new test has been applied to
# a number of patients who are known (from the old test) to have the disease or
# not. The first column of the data set (called 'prob') is the probability under
# the new test that the patient has Ebola. The second column ('ebola') is the
# result from the older test which definitively says whether the patient has ebola
# (ebola = 1) or not (ebola=0). Download the dataset, load it in and have a look
# at the first few entries to see what it looks like.
eb = pd.read_csv('ebola_test.csv')

# Q1 Write a function that returns a dict with some information about the
# DataFrame df. The keys of the dict should be 'Percentage' and 'Quartiles'. The
# value for 'Percentage' should be a single number (not a list with a number in
# it), specifying the percentage of entries with positive results for the given
# criteria, i.e. indicator variable is 1. This value should be rounded to 1
# decimal place. The value for 'Quartiles' should be a list (not a Series or
# array) with the number of observations in the 1st quartile (0-25%), 2nd
# quartile (25-50%), 3rd quartile (50-75%) and 4th quartile (75-100%) for a
# specified observation column.
# The name of the indicator and observation variable should be given to the
# function as strings.
# Be sure that your keys are exactly as specified above and that the values have
# the data type specified above.
def exercise1(df, ind_col, obs_col):
    newdict1=dict(df[ind_col].value_counts())
    obgroups=pd.qcut(df[obs_col],4)
    dict0={'Percentage':round(newdict1[1]/df.shape[0]*100,1),
          'Quartiles':list(obgroups.value_counts())}
    return dict0

# Suggested test
exercise1(eb,'ebola','prob')
# This should return
# {'Percentage': 10.4, 'Quartiles': [128, 125, 125, 121]}
# in this exact form. If it does not your function is not correct.
# Note that due to finite precision, your function may output the number
# correctly rounded to 3 decimal places, with a string of zeros followed by a
# non-zero digit, e.g. 0.070000000000000007 for the false negative rate. This is
# still a correct answer.

# In classification problems one must define a cutoff, a value for which
# observations above that value are classified as positive results (and assigned
# the value 1), and those less than or equal to the value are classified as
# negative results (and assigned the value 0). If we had a perfect classifier,
# the predicted values (0s or 1s) would match the indicator variable. A false
# positive is defined as the classifier predicting a positive result (1), when
# the actual result was negative (0). A false negative is the opposite, the
# classifier predicts a negative result (0), when the true result was positive (1).

# Q2 Write a function that takes a DataFrame, two strings specifying the names of

```

```

# indicator column and the observation column, and a cutoff value, and
# returns the rate of false positive and rate of the false negative as a dict.
# The keys of the dict should be 'False Pos' and 'False Neg' and the values
# must be rounded to 3 decimal places.
def exercise2(df, ind_col, obs_col, cutoff):
    newdict1=dict(df[ind_col][df[obs_col]>cutoff].value_counts())
    newdict2=dict(df[ind_col][df[obs_col]<=cutoff].value_counts())
    if(0 in newdict1):
        fp=round(newdict1[0]/df.shape[0],3)
    else:
        fp=0
    if(1 in newdict2):
        fn=round(newdict2[1]/df.shape[0],3)
    else:
        fn=0
    dict0={'False Pos':fp,
          'False Neg':fn}
    return dict0
# Suggested test
exercise2(eb,'ebola','prob',0.15)
# This should return
# {'False Pos': 0.126, 'False Neg': 0.07}
# or:
# {'False Neg': 0.07, 'False Pos': 0.126}
# in this exact form. If it does not your function is not correct.

# We do not know a priori what the best cutoff value is, as such, we should
# loop over a range of cutoffs to decide on the best one.

# Q3 Write a function that takes the same inputs as Q2, but cutoff will now be
# replaced with cutoff_list (an array of different cutoff values to test). The
# function should run the classification for each value in cutoff_list and
# determine which is the best cutoff value. We will define the best classifier
# as having the lowest false results (false positives plus false negatives). The
# function should return a dict with the keys 'Cutoff value', 'False Pos' and
# 'False Neg', and the values of the false positive rate and false negative rate
# should be rounded to 3 decimal places
def exercise3(df, ind_col, obs_col, cutoff_list):
    min0=1
    for cutoff in cutoff_list:
        fp=0.0
        fn=0.0
        newdict1=dict(df[ind_col][df[obs_col]>cutoff].value_counts())
        newdict2=dict(df[ind_col][df[obs_col]<=cutoff].value_counts())
        if(0 in newdict1):
            fp=newdict1[0]/df.shape[0]
        if(1 in newdict2):
            fn=newdict2[1]/df.shape[0]
        if(fp+fn<min0):
            cutoffvalue=cutoff
            min0=fp+fn
            a=fp
            b=fn
    dict0={'Cutoff value':round(cutoffvalue,3),
          'False Pos':round(a,3),
          'False Neg':round(b,3)}
    return dict0
# Suggested test
exercise3(eb,'ebola','prob',np.arange(0,1,0.01))
# This should return

```

```

# {'Cutoff value': 0.21, 'False Pos': 0.0, 'False Neg': 0.102}
# in this exact form. If it does not your function is not correct.

# A related concept to the choice of cut-offs is the Receiver Operator Characteristic
# (ROC) curve. The ROC curve aims to quantify how well a classifier beats a random
# classifier for any level of probability cut-off. An introduction can be found here:
# http://en.wikipedia.org/wiki/Receiver_operating_characteristic
# The idea is to plot the false positive rate against the true positive rate for
# every possible cut-off

# Q4 Write a function which calculates the ROC curve. The function should have
# three arguments, the DataFrame df, the name of the indicator variable ind_col
# and the name of the observation variable obs_col
# Your ROC curve function should perform the following steps
# 1) Find the unique values in the observation column
# 2) Use each of these unique values as a cutoff value and
# a) Classify all the obs as either positive or negative based on the current cutoff value
# b) Calculate the number of true positives (tp), true negatives (tn), false positives (fp) and
# Note 1: a tp is when the classification value and actual value are both 1, a tn is when they're
# Note 2: tp, fn, etc must all be vectors/Series of the same length as the vector/Series of cutoff values
# Note 4: be careful when you're at the maximum cutoff value that you can still calculate these
# 3) Create the true positive rate (tpr) as tp/(tp+fn)
# 4) Create the false positive rate (fpr) as fp/(fp+tn)
# 5) Create a DataFrame, indexed by the cutoff values (unique values of
# observation column), with columns 'True_Pos' and 'False_Pos', containing tpr
# and fpr, respectively.
# 6) Return the DataFrame sorted by index (lowest to highest)

def exercise4(df, ind_col, obs_col):
    uvalue=df[obs_col].unique()
    tpr=[]
    fpr=[]
    for cutoff in uvalue:
        newdict1=dict(df[ind_col][df[obs_col]>cutoff].value_counts())
        newdict2=dict(df[ind_col][df[obs_col]<=cutoff].value_counts())
        if(0 in newdict1):
            fp=round(newdict1[0]/df.shape[0],3)
        else:
            fp=0
        if(1 in newdict2):
            fn=round(newdict2[1]/df.shape[0],3)
        else:
            fn=0
        if(1 in newdict1):
            tp=round(newdict1[1]/df.shape[0],3)
        else:
            tp=0
        if(0 in newdict2):
            tn=round(newdict2[0]/df.shape[0],3)
        else:
            tn=0
        tpr.append(tp/(tp+fn))
        fpr.append(fp/(fp+tn))
    dict0={'cutoff values':uvalue, 'True_Pos':tpr, 'False_Pos':fpr}
    df0=pd.DataFrame(dict0, index=dict0['cutoff values'], columns=['True_Pos', 'False_Pos'])
    return df0.sort_index()

# Suggested test
Q4_ans = exercise4(eb, 'ebola', 'prob')
Q4_ans.plot(x='False_Pos', y='True_Pos')
# This should create a plot of the false positive rate vs true positive rate.
# When the false positive rate is ~0.5, the true positive rate is ~0.85

```