

```
1 # -*- coding: utf-8 -*-
2 ## Lecture 10 assessed exercises
3
4 # Packages
5 from pandas import Series, DataFrame
6 import pandas as pd
7 import numpy as np
8 import numpy.random as npr
9 import statsmodels.api as sm
10
11 # For this set of exercises we are going to use the prostate dataset and the
   diamonds
12 # dataset. Testing your functions with two different datasets should catch any
   error
13 # related to leaving the DataFrame names inside your function.
14 prostate = pd.read_csv('http://statweb.stanford.edu/~tibs/ElemStatLearn/datasets/
   prostate.data', index_col='Unnamed: 0',
15                        sep='\t')
16 diamonds = pd.read_csv('diamonds.csv')
17 # Remove categorical data and take subset of diamonds dataset
18 dataA = prostate.drop('train', axis=1)
19 dataB = dataQ1b = diamonds.drop(['cut', 'color', 'clarity'], axis=1).iloc[:100, :]
20
21
22 # Let's fit some regression models and create a stepwise AIC function
23 # As we learnt in lectures to fit a regression model, we need to create a
   DataFrame X
24 # and Series y. X should contain the standardised version of all of the
```

```
24 explanatory/
25 # exogenous variables and y should contain the standardised version of the
   response/
26 # endogenous variable. To fit the intercept, X must have an additional column of
   ones.
27
28 # Q1 Write a function to create X and y for a given DataFrame df. The function
   inputs are
29 # the DataFrame df and the label of the response/endogenous variable res_col. The
   function
30 # should return two objects, X and y (in that order), where X and y are both
   standardised
31 # and the column of ones is the first column of X.
32 # (You may assume that none of the variables are categorical)
33 def exercise1(df, res_col):
34     df = (df - df.mean()) / df.std()
35     y = df[res_col]
36     X = df.drop(res_col, axis=1)
37     X = (X - X.mean()) / X.std()
38     X.insert(0, 'intercept', np.ones(len(X)))
39     return X, y
40
41
42 # Suggested tests
43 XA, yA = exercise1(dataA, 'lpsa')
44 XB, yB = exercise1(dataB, 'price')
45
46
```

```
47 # Things to check to ensure code is working correctly
48 # - XA and XB have the same number of rows and columns as dataA and dataB,
    respectively
49 # - the first column of XA and XB is entirely ones
50 # - yA and yB are Series with the same number of rows as dataA and dataB,
    respectively
51 # - the mean of each variable in XA, XB, yA and yB (apart from the intercept
    column)
52 # is close to zero ( $\sim 10^{-16}$ )
53 # - the std of each variable in XA, XB, yA and yB (apart from the intercept column
    )
54 # is 1
55
56 # Q2 Write a function that takes X and y as inputs and fits a linear regression
    model.
57 # The function should return the rsquared value rounded to 4 decimal places
58 def exercise2(X, y):
59     mod = sm.OLS(y, X)
60     res = mod.fit()
61     return round(res.rsquared, 4)
62
63
64 # Suggested tests
65 # Remember we can unpack a tuple to use as a set of inputs to a function. Here we
    unpack
66 # the tuple (X,y) returned by exercise1 to use as an input for exercise2
67 print(exercise2(*exercise1(dataA, 'lpsa'))))
68 # Should give 0.6634
```

```
69 print(exercise2(*exercise1(dataB, 'price')))  
70 # Should give 0.9426  
71  
72  
73 # AIC is the Akaike information criterion. It's designed to penalise models with  
74 # lots of explanatory variables so that we pick models which fit the data well but  
75 # aren't too complicated. In general, if you have two models fitted to the same  
76 # data,  
77 # the model with the lowest AIC is preferable. The AIC is given as part of the  
78 # model  
79 # summary with OLS  
80 # The steps to run a forward selection AIC regression are:  
81 # 1. Run a linear regression with just the intercept column. Get the AIC.  
82 # 2. Add in the explanatory variables individually, run a linear regression for  
83 # each one and determine  
84 # how much they decreases the AIC  
85 # 3. Find the variable with the biggest decrease in AIC and include it in your  
86 # linear model  
87 # 4. Repeat step 2-3 with this new linear model and remaining explanatory  
88 # variables  
89 # 5. Repeat this process until none of the remaining explanatory variables reduce  
90 # the AIC  
91 # The explanatory variables that have been included up to the stopping point are  
92 # considered the  
93 # variables that produce a good fit without overcomplicating the model.  
94  
95 # Q3 Write a function that performs the AIC algorithm for a given DataFrame X and
```

```
89 Series y.
90 # The function should return the names of the columns used for the model that
   gives the lowest AIC
91 # This question is worth 2 marks
92 def exercise3(X, y):
93     mod = sm.OLS(y, X['intercept'])
94     res = mod.fit()
95     aic = res.aic
96     ex_v = ['intercept']
97     col_name = [n for n in X.columns]
98     col_name.remove('intercept')
99     while True:
100         tem_aic = []
101         for i in range(len(col_name)):
102             tem_v = ex_v.copy()
103             tem_v.append(col_name[i])
104             mod = sm.OLS(y, X[tem_v])
105             tem_aic.append(mod.fit().aic)
106         if min(tem_aic) < aic:
107             aic = min(tem_aic)
108             ex_v.append(col_name[np.argmin(tem_aic)])
109             col_name.remove(col_name[np.argmin(tem_aic)])
110         else:
111             break
112     return ex_v
113
114
115 # Suggested tests
```

```
116 print(exercise3(*exercise1(dataA, 'lpsa')))  
117 # Should give ['intercept', 'lcavol', 'lweight', 'svi', 'lbph', 'age']  
118 print(exercise3(*exercise1(dataB, 'price')))  
119 # Should give ['intercept', 'carat', 'z', 'x', 'y', 'table']  
120
```