

```

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 from scipy import stats
11
12 # For this set of exercises we are going to use the prostate dataset and the
13 # diamonds dataset. Testing your functions with two different datasets should catch any error
14 # related to leaving the DataFrame names inside your function.
15 prostate = pd.read_csv(
16     'http://statweb.stanford.edu/~tibs/ElemStatLearn/datasets/prostate.data',
17     index_col='Unnamed: 0', sep='\t')
18 diamonds = pd.read_csv('./Diamonds.csv')
19 # Remove categorical data and take subset of diamonds dataset
20 dataA = prostate.drop('train', axis=1)
21 dataB = dataQ1b = diamonds.drop(
22     ['cut', 'color', 'clarity'], axis=1).iloc[:100, :]
23
24 # Let's fit some regression models and create a stepwise AIC function
25 # As we learnt in lectures to fit a regression model, we need to create a DataFrame
26 # X
27 # and Series y. X should contain the standardised version of all of the explanatory/
28 # exogenous variables and y should contain the standardised version of the response/
29 # endogenous variable. To fit the intercept, X must have an additional column of
30 # ones.
31
32 # Q1 Write a function to create X and y for a given DataFrame df. The function
33 # inputs are
34 # the DataFrame df and the label of the response/endogenous variable res_col. The
35 # function
36 # should return two objects, X and y (in that order), where X and why are both
37 # standardised
38 # and the column of ones is the first column of X.
39 # (You may assume that none of the variables are categorical)
40
41 def exercise1(df, res_col):
42     # Reset the index of DataFrame and make it start from 0.
43     df = df.reset_index(drop=True)
44     # Data standardization
45     df = (df-df.mean())/df.std()
46     y = df.pop(res_col)
47     n = df.shape[0]
48     # create a column of ones and change it into dataframe
49     X1 = np.ones((n, 1))
50     X1 = pd.DataFrame(data=X1, columns=["intercept"], index=range(n))
51     df.reindex_like(X1)
52     # combine two dataframes
53     X = pd.concat([X1, df], axis=1)
54     return X, y
55
56 # Suggested tests
57 XA, yA = exercise1(dataA, 'lpsa')
58 XB, yB = exercise1(dataB, 'price')

```

```

55 # Things to check to ensure code is working correctly
56 # - XA and XB have the same number of rows and columns as dataA and dataB,
    respectively
57 # - the first column of XA and XB is entirely ones
58 # - yA and yB are Series with the same number of rows as dataA and dataB,
    respectively
59 # - the mean of each variable in XA, XB, yA and yB (apart from the intercept column)
60 # is close to zero ( $\sim 10^{-16}$ )
61 # - the std of each variable in XA, XB, yA and yB (apart from the intercept column)
62 # is 1
63
64 # Q2 Write a function that takes X and y as inputs and fits a linear regression
    model.
65 # The function should return the rsquared value rounded to 4 decimal places
66
67
68 def exercise2(X, y):
69     # create a linear regression model then get the value of R-s
70     mod = sm.OLS(y, X)
71     # then get the value of R-squared
72     r2 = mod.fit().rsquared
73     r2 = round(r2, 4)
74     return r2
75
76
77 # Suggested tests
78 # Remember we can unpack a tuple to use as a set of inputs to a function. Here we
    unpack
79 # the tuple (X,y) returned by exercise1 to use as an input for exercise2
80 print(exercise2(*exercise1(dataA, 'lpsa')))
81 # Should give 0.6634
82 print(exercise2(*exercise1(dataB, 'price')))
83 # Should give 0.9426
84
85
86 # AIC is the Akaike information criterion. It's designed to penalise models with
87 # lots of explanatory variables so that we pick models which fit the data well but
88 # aren't too complicated. In general, if you have two models fitted to the same
    data,
89 # the model with the lowest AIC is preferable. The AIC is given as part of the model
90 # summary with OLS
91
92 # The steps to run a forward selection AIC regression are:
93 # 1. Run a linear regression with just the intercept column. Get the AIC.
94 # 2. Add in the explanatory variables individually, run a linear regression for each
    one and determine
95 #     how much they decreases the AIC
96 # 3. Find the variable with the biggest decrease in AIC and include it in your
    linear model
97 # 4. Repeat step 2-3 with this new linear model and remaining explanatory variables
98 # 5. Repeat this process until none of the remaining explanatory variables reduce
    the AIC
99 # The explanatory variables that have been included up to the stopping point are
    considered the
100 # variables that produce a good fit without overcomplicating the model.
101
102 # Q3 Write a function that performs the AIC algorithm for a given DataFrame X and
    Series y.
103 # The function should return the names of the columns used for the model that gives
    the lowest AIC
104 # This question is worth 2 marks

```

```

105 def exercise3(X, y):
106     # label: a vector that helps to loop, including the column names of X.
107     # col_select: a list used to store the column names that have the least aic
values when regressing.
108     # X1: the selected values of variables in X, including the intercept column at
first.
109     # X2: the the rest values of variables in X.
110     # aic_flag: the minimum AIC values of models, which is initialed as regressing
with just the intercept column.
111     label = X.columns
112     col_select = list()
113     col_select.append(label[0])
114     label = np.delete(label, 0, 0)
115     y = y.values
116     X1 = X.pop("intercept").values
117     X2 = X.values
118     aic_flag = sm.OLS(y, X1.T).fit().aic
119
120     while(any(label)):
121         aic_select = list()
122         # get each AIC value of variable in X2, when creating linear regression with
X1.
123         for i in range(X2.shape[1]):
124             X3 = np.vstack((X1, X2[:, i]))
125             X3 = X3.T
126             aic_val = sm.OLS(y, X3).fit().aic
127             aic_select.append(aic_val)
128             # searching for the least AIC value
129             index = np.argmin(aic_select)
130             # If the minimum AIC value is less than current one
131             # reset the minimum AIC value
132             # select this variable, and add its values to X1.
133             if min(aic_select) < aic_flag:
134                 aic_flag = min(aic_select)
135                 col_select.append(label[index])
136                 X1 = np.vstack((X1, X2[:, index]))
137             # Deleting the variable that has evaluated. When the label becomes none, the
loop will stop
138             X2 = np.delete(X2, index, 1)
139             label = np.delete(label, index, 0)
140         return col_select
141
142
143 # Suggested tests
144 print(exercise3(*exercise1(dataA, 'lpsa')))
145 # Should give ['intercept', 'lcavol', 'lweight', 'svi', 'lbph', 'age']
146 print(exercise3(*exercise1(dataB, 'price')))
147 # Should give ['intercept', 'carat', 'z', 'x', 'y', 'table']
148
149

```