



Join GitHub today

GitHub is home to over 36 million developers working together to host and review code, manage projects, and build software together.

Sign up

Dismiss

Branch: master ▾

Introduction-To-Data-Science-In-Python / Assignment+2.py

Find file

Copy path



Satya Assignment 2

88888f4 on Aug 1, 2018

0 contributors

204 lines (148 sloc) | 7.48 KB

Raw

Blame

History



```
1
2 # coding: utf-8
3
4 # ---
5 #
6 # You are currently looking at **version 1.2** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks, please visit https://mybinder.org
7 #
8 # ---
9
10 # Assignment 2 - Pandas Introduction
11 # All questions are weighted the same in this assignment.
12 # ## Part 1
13 # The following code loads the olympics dataset (olympics.csv), which was derived from the Wikipedia entry on [All Time Olympic Games Medals](https://en.wikipedia.org/wiki/List_of_all-time_olympic_medalists).
14 #
15 # The columns are organized as # of Summer games, Summer medals, # of Winter games, Winter medals, total # number of games, total # of medals.
16
17 # In[5]:
18
19 import pandas as pd
20
21 df = pd.read_csv('olympics.csv', index_col=0, skiprows=1)
22
23 for col in df.columns:
24     if col[:2]=='01':
25         df.rename(columns={col:'Gold'+col[4:]}, inplace=True)
26     if col[:2]=='02':
27         df.rename(columns={col:'Silver'+col[4:]}, inplace=True)
28     if col[:2]=='03':
29         df.rename(columns={col:'Bronze'+col[4:]}, inplace=True)
30     if col[:1]=='#':
31         df.rename(columns={col:'#'+col[1:]}, inplace=True)
32
33 names_ids = df.index.str.split('\s\(') # split the index by '('
34 #names_ids
35 df.index = names_ids.str[0] # the [0] element is the country name (new index)
36 df['ID'] = names_ids.str[1].str[:3] # the [1] element is the abbreviation or ID (take first 3 characters from that)
37
38 df = df.drop('Totals')
39 df.head(10)
40
41
42 # ### Question 0 (Example)
43 #
44 # What is the first country in df?
45 #
46 # *This function should return a Series.*
47
48 # In[6]:
49
50 # You should write your whole answer within the function provided. The autograder will call
51 # this function and compare the return value against the correct solution value
52 def answer_zero():
53     # This function returns the row for Afghanistan, which is a Series object. The assignment
54     # question description will tell you the general format the autograder is expecting
55     return df.iloc[0]
56
57 # You can examine what your function returns by calling it in the cell. If you have questions
58 # about the assignment formats, check out the discussion forums for any FAQs
59 answer_zero()
60
61
```

```

62 # ### Question 1
63 # Which country has won the most gold medals in summer games?
64 #
65 # *This function should return a single string value.*
66
67 # In[7]:
68
69 def answer_one():
70     max_gold = df['Gold'].max()
71     df1 = df[df['Gold']==max_gold]
72     return str(df1.iloc[0].name)
73
74
75 # ### Question 2
76 # Which country had the biggest difference between their summer and winter gold medal counts?
77 #
78 # *This function should return a single string value.*
79
80 # In[8]:
81
82 def answer_two():
83     #import numpy as np
84     diff = df['Gold'] - df['Gold.1']
85     diff=diff.abs()
86     res=diff.max()
87     diff=diff[diff==res]
88     return str(diff.index.values[0])
89
90
91 # ### Question 3
92 # Which country has the biggest difference between their summer gold medal counts and winter gold medal counts relative to their total gold
93 #
94 #  $\frac{\text{Summer-Gold} - \text{Winter-Gold}}{\text{Total-Gold}}$ 
95 #
96 # Only include countries that have won at least 1 gold in both summer and winter.
97 #
98 # *This function should return a single string value.*
99
100 # In[96]:
101
102 def answer_three():
103     res=df[(df['Gold']>0) & (df['Gold.1']>0)]
104     res = (res['Gold']-res['Gold.1'])/(res['Gold'] + res['Gold.1'] + res['Gold.2'])
105     res=res[res==res.max()]
106     return res.index.values[0]
107
108 answer_three()
109
110
111 # ### Question 4
112 # Write a function that creates a Series called "Points" which is a weighted value where each gold medal ('Gold.2') counts for 3 points, si
113 #
114 # *This function should return a Series named 'Points' of length 146*
115
116 # In[10]:
117
118 def answer_four():
119     res = df['Gold.2']*3 + df['Silver.2']*2 + df['Bronze.2']
120     res.name = 'Points'
121     return res
122
123
124 # ## Part 2
125 # For the next set of questions, we will be using census data from the [United States Census Bureau](http://www.census.gov). Counties are p
126 #
127 # The census dataset (census.csv) should be loaded as census_df. Answer questions using this as appropriate.
128 #
129 # ### Question 5
130 # Which state has the most counties in it? (hint: consider the sumlevel key carefully! You'll need this for future questions too...)
131 #
132 # *This function should return a single string value.*
133
134 # In[11]:
135
136 census_df = pd.read_csv('census.csv')
137 census_df.head()
138
139
140 # In[12]:
141
142 def answer_five():
143     res=census_df[census_df['SUMLEV']==50]
144     res=res.groupby('STNAME').count()['SUMLEV']
145     return res[res==res.max()].index.values[0]
146
147
148 # ### Question 6
149 # **Only looking at the three most populous counties for each state**, what are the three most populous states (in order of highest populat
150 #
151 # *This function should return a list of string values *

```

```

152
153 # In[39]:
154
155 def answer_six():
156     res=census_df[['SUMLEV', 'STNAME', 'CENSUS2010POP']]
157     res=res[res['SUMLEV']==50]
158     res=res.sort(columns=['STNAME', 'CENSUS2010POP'], ascending=[True,False])
159     res=res.groupby('STNAME').head(3)
160     res=res.groupby('STNAME').sum()
161     res=res.nlargest(3,'CENSUS2010POP')
162     return list(res.index.values)
163
164
165 # ### Question 7
166 # Which county has had the largest absolute change in population within the period 2010-2015? (Hint: population values are stored in column
167 #
168 # e.g. If County Population in the 5 year period is 100, 120, 80, 105, 100, 130, then its largest change in the period would be |130-80| =
169 #
170 # *This function should return a single string value.*
171
172 # In[98]:
173
174 def answer_seven():
175     import numpy as np
176     res=census_df[census_df['SUMLEV']==50]
177     res=res[['CTYNAME', 'POPESTIMATE2010', 'POPESTIMATE2011', 'POPESTIMATE2012', 'POPESTIMATE2013', 'POPESTIMATE2014', 'POPESTIMATE2015']]
178     max_p=map(max,res['POPESTIMATE2010'],res['POPESTIMATE2011'],res['POPESTIMATE2012'],res['POPESTIMATE2013'],res['POPESTIMATE2014'],res['P
179     max_p=np.array(list(max_p))
180     min_p=map(min,res['POPESTIMATE2010'],res['POPESTIMATE2011'],res['POPESTIMATE2012'],res['POPESTIMATE2013'],res['POPESTIMATE2014'],res['P
181     min_p=np.array(list(min_p))
182     chg=max_p-min_p
183     return res.iloc[chg.argmax()].values[0]
184
185
186 # ### Question 8
187 # In this datafile, the United States is broken up into four regions using the "REGION" column.
188 #
189 # Create a query that finds the counties that belong to regions 1 or 2, whose name starts with 'Washington', and whose POPESTIMATE2015 was
190 #
191 # *This function should return a 5x2 DataFrame with the columns = ['STNAME', 'CTYNAME'] and the same index ID as the census_df (sorted asce
192
193 # In[94]:
194
195 def answer_eight():
196     res=census_df[(census_df['SUMLEV']==50) & (census_df['REGION'] <= 2) & (census_df['CTYNAME'].str.match('^Washington')) & (census_df['PO
197     return res[['STNAME', 'CTYNAME']].sort_index()
198
199
200 # In[ ]:
201
202
203

```

