

# Assignment 2 - Pandas Introduction

All questions are weighted the same in this assignment.

## Part 1

The following code loads the olympics dataset (olympics.csv), which was derived from the Wikipedia entry on [All Time Olympic Games Medals](#), and does some basic data cleaning.

The columns are organized as # of Summer games, Summer medals, # of Winter games, Winter medals, total # number of games, total # of medals. Use this dataset to answer the questions below.

```
In [12]: import pandas as pd

df = pd.read_csv('olympics.csv', index_col=0, skiprows=1)

for col in df.columns:
    if col[:2]=='01':
        df.rename(columns={col:'Gold'+col[4:]}, inplace=True)
    if col[:2]=='02':
        df.rename(columns={col:'Silver'+col[4:]}, inplace=True)
    if col[:2]=='03':
        df.rename(columns={col:'Bronze'+col[4:]}, inplace=True)
    if col[:1]=='W':
        df.rename(columns={col:'#'+col[1:]}, inplace=True)

names_ids = df.index.str.split('\s\(') # split the index by '('

df.index = names_ids.str[0] # the [0] element is the country name (new index)
df['ID'] = names_ids.str[1].str[:3] # the [1] element is the abbreviation or ID (take first 3 characters from that)

df = df.drop('Totals')
df.head()
```

Out[12]:

	# Summer	Gold	Silver	Bronze	Total	# Winter	Gold.1	Silver.1	Bronze.1	Total.1	# Games	Gold.2	Silve
Afghanistan	13	0	0	2	2	0	0	0	0	0	13	0	0
Algeria	12	5	2	8	15	3	0	0	0	0	15	5	2
Argentina	23	18	24	28	70	18	0	0	0	0	41	18	24
Armenia	5	1	2	9	12	6	0	0	0	0	11	1	2
Australasia	2	3	4	5	12	0	0	0	0	0	2	3	4

### Question 0 (Example)

What is the first country in df?

This function should return a Series.

```
In [17]: # You should write your whole answer within the function provided. The autograder will call
# this function and compare the return value against the correct solution value
def answer_zero():

    # This function returns the row for Afghanistan, which is a Series object. The assignment
    # question description will tell you the general format the autograder is expecting
    return df.iloc[0]

# You can examine what your function returns by calling it in the cell. If you have questions
# about the assignment formats, check out the discussion forums for any FAQs
answer_zero()
```

Out[17]:

# Summer	13
Gold	0
Silver	0
Bronze	2
Total	2
# Winter	0
Gold.1	0
Silver.1	0
Bronze.1	0
Total.1	0
# Games	13
Gold.2	0
Silver.2	0
Bronze.2	2
Combined total	2
ID	AFG
Name: Afghanistan, dtype: object	

### Question 1

Which country has won the most gold medals in summer games?

This function should return a single string value.

```
In [19]: def answer_one():
maxG=df[df['Gold']==max(df['Gold'])].index.tolist()[0]
return maxG
answer_one()
```

Out[19]: 'United States'

### Question 2

Which country had the biggest difference between their summer and winter gold medal counts?

This function should return a single string value.

```
In [20]: def answer_two():
return df[(df['Gold']-df['Gold.1'])==max(df['Gold']-df['Gold.1'])].index.tolist()[0]
answer_two()
```

Out[20]: 'United States'

### Question 3

Which country has the biggest difference between their summer gold medal counts and winter gold medal counts relative to their total gold medal count?

$$\frac{\text{Summer Gold} - \text{Winter Gold}}{\text{Total Gold}}$$

Only include countries that have won at least 1 gold in both summer and winter.

This function should return a single string value.

```
In [19]: def answer_three():
return ((df['Gold'].where(df['Gold']>0)-df['Gold.1'].where(df['Gold.1']>0))/df['Gold.2']).argmax()
answer_three()#was getting "Array conditional must be same shape as self" when using df['Gold'].count
```

Out[19]: 'Bulgaria'

### Question 4

Write a function that creates a Series called "Points" which is a weighted value where each gold medal (Gold.2) counts for 3 points, silver medals (Silver.2) for 2 points, and bronze medals (Bronze.2) for 1 point. The function should return only the column (a Series object) which you created, with the country names as indices.

This function should return a Series named Points of length 146

```
In [26]: def answer_four():
Points=pd.Series (3*df['Gold.2']+2*df['Silver.2']+df['Bronze.2'],index=df.index)
return Points
answer_four()
```

Out[26]:

Afghanistan	2
Algeria	12
Argentina	76
Armenia	13
Australasia	13
Australia	496
Austria	397
Azerbaijan	25
Bahamas	9
Bahrain	1
Barbados	1
Belarus	106
Belgium	163
Bermuda	1
Bohemia	5
Botswana	2
Brazil	115
British West Indies	2
Bulgaria	256
Burundi	0
Cameroon	3
Canada	545
Chile	18
China	493
Colombia	23
Costa Rica	4
Ivory Coast	2
Croatia	41
Cuba	204
Cyprus	2
...	
Spain	155
Sri Lanka	4
Sudan	2
Suriname	1
Sweden	688
Switzerland	389
Syria	3
Chinese Taipei	26
Tajikistan	4
Tanzania	4
Thailand	23
Togo	1
Tonga	2
Trinidad and Tobago	21
Tunisia	10
Turkey	74
Uganda	8
Ukraine	117
United Arab Emirates	0
United States	2564
Uruguay	10
Uzbekistan	21
Venezuela	12
Vietnam	4
Virgin Islands	2
Yugoslavia	93
Independent Olympic Participants	4
Zambia	3
Zimbabwe	9
Mixed team	14
dtype: int64	

## Part 2

For the next set of questions, we will be using census data from the [United States Census Bureau](#). Counties are political and geographic subdivisions of states in the United States. This dataset contains population data for counties and states in the US from 2010 to 2015. [See this document](#) for a description of the variable names.

The census dataset (census.csv) should be loaded as census\_df. Answer questions using this as appropriate.

### Question 5

Which state has the most counties in it? (hint: consider the sumlevel key carefully! You'll need this for future questions too...)

This function should return a single string value.

```
In [1]: import pandas as pd
census_df = pd.read_csv('census.csv')
census_df.head()
```

Out[1]:

	SUMLEV	REGION	DIVISION	STATE	COUNTY	STNAME	CTYNAME	CENSUS2010POP	ESTIMATESBASE2010	P
0	40	3	6	1	0	Alabama	Alabama	4779736	4780127	4
1	50	3	6	1	1	Alabama	Autauga County	54571	54571	5
2	50	3	6	1	3	Alabama	Baldwin County	182265	182265	1
3	50	3	6	1	5	Alabama	Barbour County	27457	27457	2
4	50	3	6	1	7	Alabama	Bibb County	22915	22919	2

5 rows x 100 columns

```
In [9]: def answer_five():
#id=df['COUNTY'].argmax()
counties_df = census_df[census_df['SUMLEV'] == 50]

return counties_df.groupby('STNAME').count()['CTYNAME'].idxmax()
answer_five()
```

Out[9]: 'West Virginia'

### Question 6

Only looking at the three most populous counties for each state, what are the three most populous states (in order of highest population to lowest population)? Use CENSUS2010POP.

This function should return a list of string values.

```
In [26]: def answer_six():
counties_df = census_df[census_df['SUMLEV'] == 50]
top_three=counties_df.sort_values(by='CENSUS2010POP', ascending=False).groupby('STNAME').head(3)
return top_three.groupby('STNAME').sum().sort_values(by='CENSUS2010POP', ascending=False).head(3)
).index.tolist()
answer_six()
```

Out[26]: ['California', 'Texas', 'Illinois']

### Question 7

Which county has had the largest absolute change in population within the period 2010-2015? (Hint: population values are stored in columns POPESTIMATE2010 through POPESTIMATE2015, you need to consider all six columns.)

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| = 50.

This function should return a single string value.

```
In [8]: def answer_seven():
import numpy as np
counties_df = census_df[census_df['SUMLEV'] == 50]
counties_df =counties_df.set_index('CTYNAME')
columns=['POPESTIMATE2010','POPESTIMATE2011','POPESTIMATE2012','POPESTIMATE2013','POPESTIMATE2014','POPESTIMATE2015']
minct=counties_df[columns].min(axis=1)
maxct=counties_df[columns].max(axis=1)
diffct=maxct-minct
counties_df['DIFF']=diffct
ct=counties_df['DIFF'].idxmax()
return ct
answer_seven()
```

Out[8]: 'Harris County'

### Question 8

In this datafile, the United States is broken up into four regions using the "REGION" column.

Create a query that finds the counties that belong to regions 1 or 2, whose name starts with 'Washington', and whose POPESTIMATE2015 was greater than their POPESTIMATE 2014.

This function should return a 5x2 DataFrame with the columns = ['STNAME', 'CTYNAME'] and the same index ID as the census\_df (sorted ascending by index).

```
In [61]: def answer_eight():
counties_df = census_df[census_df['SUMLEV'] == 50]
new_df=counties_df[(counties_df['REGION'] ==1 ) | (counties_df['REGION'] ==2)]&(counties_df['CTY NAME'] == 'Washington County') &(counties_df['POPESTIMATE2015']>counties_df['POPESTIMATE2014'])[['STNAME','CTYNAME']]
return new_df
answer_eight()
```

Out[61]:

	STNAME	CTYNAME
896	Iowa	Washington County
1419	Minnesota	Washington County