

Unit Testing

While we will not cover the [unit testing library](https://docs.python.org/3/library/unittest.html) (<https://docs.python.org/3/library/unittest.html>) that python has, we wanted to introduce you to a simple way that you can test your code.

Unit testing is important because it is the only way you can be sure that your code is doing what you think it is doing.

Remember, just because there are no errors does not mean your code is correct.

```
In [1]: import numpy as np
import pandas as pd
import matplotlib as plt
pd.set_option('display.max_columns', 100) # Show all columns when looking at dataframe
```

```
In [2]: # Download NHANES 2015-2016 data
df = pd.read_csv("nhanes_2015_2016.csv")
df.index = range(1,df.shape[0]+1)
```

```
In [3]: df.head()
```

Out[3]:

	SEQN	ALQ101	ALQ110	ALQ130	SMQ020	RIAGENDR	RIDAGEYR	RIDRETH1	DMDCITZN	DMDEDUC2	DMDMARTL	DMDHHSIZ	WTI
1	83732	1.0	NaN	1.0	1	1	62	3	1.0	5.0	1.0	2	1346
2	83733	1.0	NaN	6.0	1	1	53	3	2.0	3.0	3.0	1	2432
3	83734	1.0	NaN	NaN	1	1	78	3	1.0	3.0	1.0	2	1240
4	83735	2.0	1.0	1.0	2	2	56	3	1.0	5.0	6.0	1	1027
5	83736	2.0	1.0	1.0	2	2	42	4	1.0	4.0	3.0	5	1762

Goal

We want to find the mean of first 100 rows of 'BPXSY1' when 'RIDAGEYR' > 60

```
In [14]: # One possible way of doing this is:
pd.Series.mean(df[df.RIDAGEYR > 60].loc[range(0,100), 'BPXSY1'])

# Current version of python will include this warning, older versions will not
# Done by Pandas becoz numpy will take more than one line to do the same
```

Out[14]: 139.57142857142858

```
In [15]: # test our code on only ten rows so we can easily check
test = pd.DataFrame({'col1': np.repeat([3,1],5), 'col2': range(3,13)}, index=range(1,11))
test
```

Out[15]:

	col1	col2
1	3	3
2	3	4
3	3	5
4	3	6
5	3	7
6	1	8
7	1	9
8	1	10
9	1	11
10	1	12

```
In [16]: # pd.Series.mean(df[df.RIDAGEYR > 60].loc[range(0,5), 'BPXSY1'])
# should return 5

pd.Series.mean(test[test.col1 > 2].loc[range(0,5), 'col2'])
```

Out[16]: 4.5

What went wrong?

```
In [17]: test[test.col1 > 2].loc[range(0,5), 'col2']
# 0 is not in the row index labels because the second row's value is < 2. For now, pandas defaults to filling
this
# with NaN
```

```
Out[17]: 0    NaN
1    3.0
2    4.0
3    5.0
4    6.0
Name: col2, dtype: float64
```

```
In [18]: # Using the .iloc method instead, we are correctly choosing the first 5 rows, regardless of their row labels
test[test.col1 >2].iloc[range(0,5), 1]
```

```
Out[18]: 1    3
2    4
3    5
4    6
5    7
Name: col2, dtype: int64
```

```
In [19]: pd.Series.mean(test[test.col1 >2].iloc[range(0,5), 1])
```

```
Out[19]: 5.0
```

```
In [20]: # We can compare what our real dataframe looks like with the incorrect and correct methods
df[df.RIDAGEYR > 60].loc[range(0,5), :] # Filled with NaN whenever a row label does not meet the condition
```

```
Out[20]:
```

	SEQN	ALQ101	ALQ110	ALQ130	SMQ020	RIAGENDR	RIDAGEYR	RIDRETH1	DMDCITZN	DMDEDUC2	DMDMARTL	DMDHHSIZ	WT
0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	Na
1	83732.0	1.0	NaN	1.0	1.0	1.0	62.0	3.0	1.0	5.0	1.0	2.0	13
2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	Na
3	83734.0	1.0	NaN	NaN	1.0	1.0	78.0	3.0	1.0	3.0	1.0	2.0	12
4	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	Na

```
In [21]: df[df.RIDAGEYR > 60].iloc[range(0,5), :] # Correct picks the first five rows such that 'RIDAGEYR' > 60
```

```
Out[21]:
```

	SEQN	ALQ101	ALQ110	ALQ130	SMQ020	RIAGENDR	RIDAGEYR	RIDRETH1	DMDCITZN	DMDEDUC2	DMDMARTL	DMDHHSIZ	WT
1	83732	1.0	NaN	1.0	1	1	62	3	1.0	5.0	1.0	2	13
3	83734	1.0	NaN	NaN	1	1	78	3	1.0	3.0	1.0	2	12
6	83737	2.0	2.0	NaN	2	2	72	1	2.0	2.0	4.0	5	11
14	83754	2.0	1.0	1.0	2	2	67	2	1.0	5.0	1.0	7	10
15	83755	1.0	NaN	3.0	2	1	67	4	1.0	5.0	2.0	1	14

```
In [22]: # Applying the correct method to the original question about BPXSY1
print(pd.Series.mean(df[df.RIDAGEYR > 60].iloc[range(0,100), 16]))
```

```
# Another way to reference the BPXSY1 variable
print(pd.Series.mean(df[df.RIDAGEYR > 60].iloc[range(0,100), df.columns.get_loc('BPXSY1')]))
```

```
136.29166666666666
136.29166666666666
```