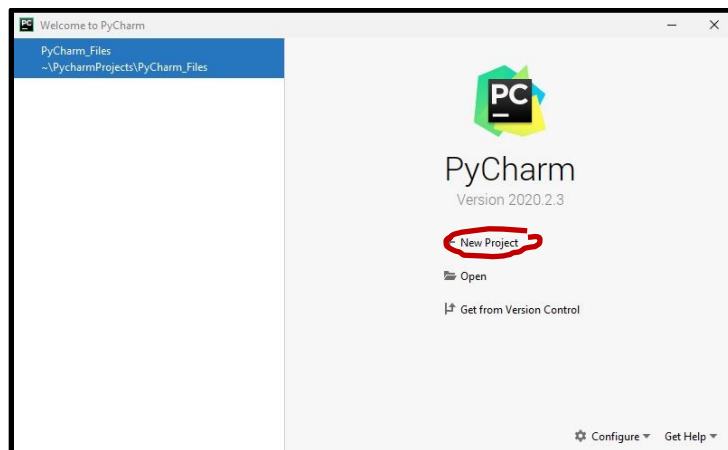# Data Wrangling

## ➕ What constitutes Data Wrangling?

Data wrangling involves processing data in various formats like - merging, grouping, concatenating etc. for the purpose of analyzing or getting data ready to be used with another data set. Data wrangling is different from data profiling.  Data profiling is used for data familiarity and inspection. Data extraction, transformation, and loading (ETL) is formalized data wrangling to collect data from operational data stores and other external sources, wrangle the data into standardized formats, and write the data into a permanent data warehouse archive.  Python has built-in features to help in data wrangling methods for various data sets.  For more on data wrangling not in this tutorial, refer to the URL:
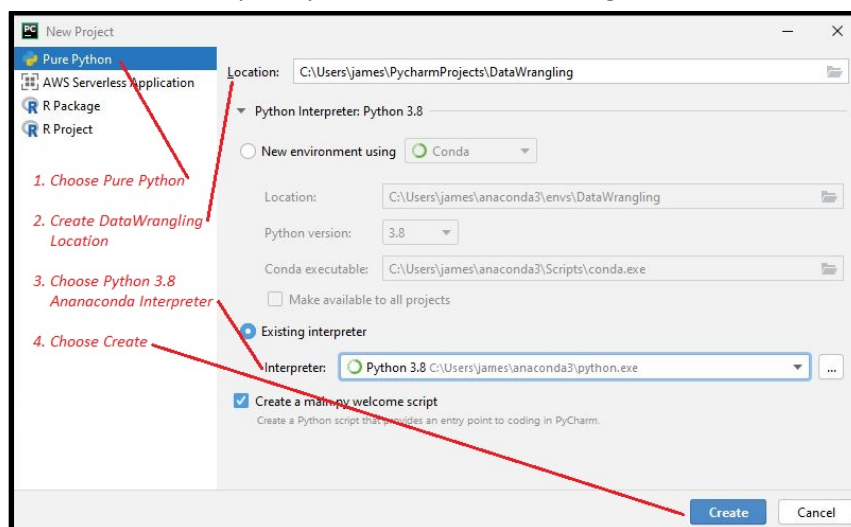
https://www.tutorialspoint.com/python_data_science/python_data_wrangling.htm
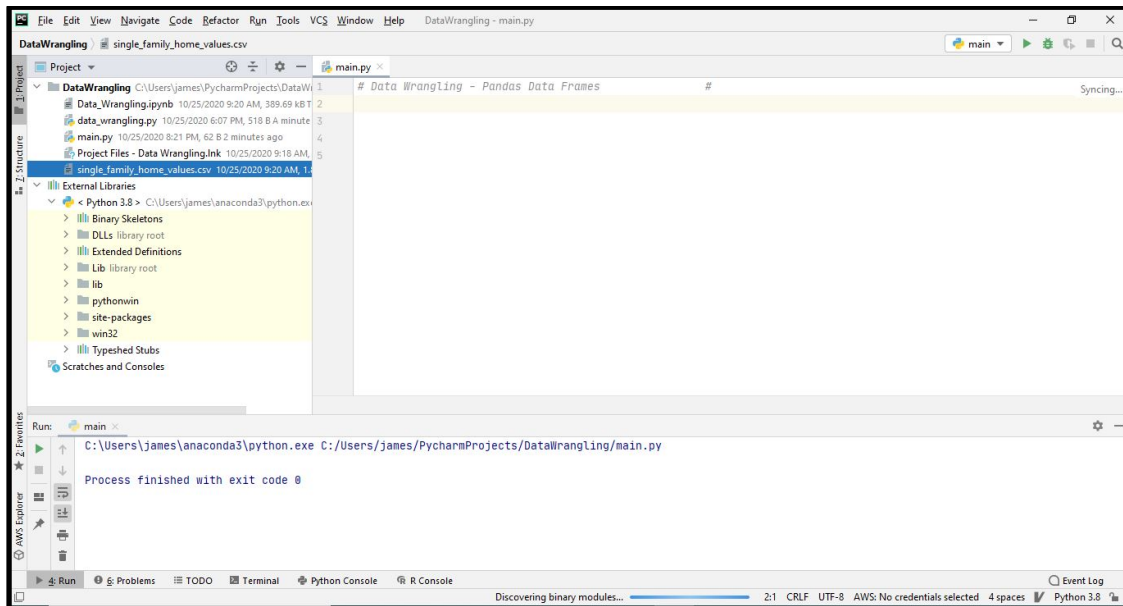
## ➕ Create a Data Wrangling Project in PyCharm

0)  Open PyCharm and choose a New Project.



1)  Choose a Pure Python project,

2)  Change your project location to create a DataWrangling folder

3)  Be sure to choose the Python interpreter as Python 3.8 associated with Anaconda.

4)  Choose Create to open PyCharm with these settings.

5)  *Download* the single_family_home_values.csv file from the Week 09 Canvas Module to the DataWrangling directory created in the last step.



## 🔸 Reading a CSV file into a Pandas DataFrame:

Pandas is an open-source Python Library used for high-performance data manipulation and data analysis using its powerful data structures. Pandas use Series and DataFrame on top of Numpy arrary. In this tutorial, we will look at DataFrame or df. Python with pandas is in use in a variety of academic and commercial domains, including Finance, Economics, Statistics, Advertising, Web Analytics, and more. Using Pandas, we can accomplish five typical steps in the processing and analysis of data, regardless of the origin of data — load, organize, manipulate, model, and analyze the data. For more specific details about Pandas not covered in this tutorial, please refer to the following URL: https://www.tutorialspoint.com/python_data_science/python_pandas.htm

1)  Anaconda will load the libraries in PyCharm for you. To access the three libraries we are using for Tutorial Seven (e.g., pandas, seaborn, and matplotlib.pyplot), enter the following code into Main.py in the DataWrangling project directory of PyCharm:

```
import pandas as pd
import seaborn as sn
import matplotlib.pyplot as plt
```

2)  To format the run window to handle the size of the csv data file we will be wrangling in Tutorial Seven, enter the following Python code:
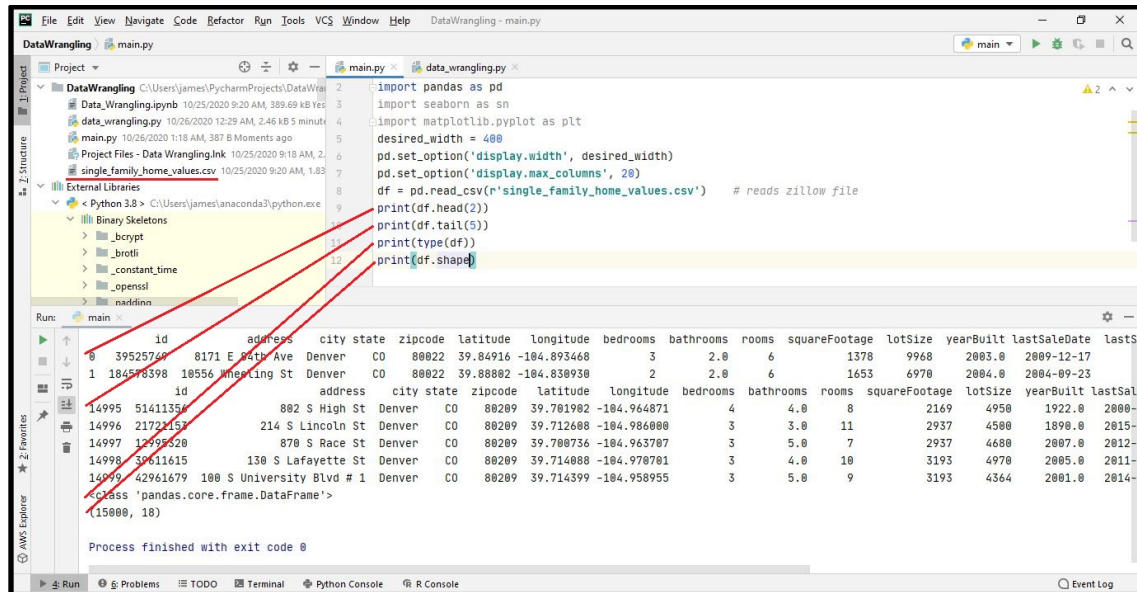
```
desired_width = 40
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
```

3)  To read the single_family_home_values.csv file into a pandas dataframe, type the following Python code into main.py via PyCharm in the line below the previous code above.

```
df = pd.read_csv(r'single_family_home_values.csv')    # reads Zillow file
```
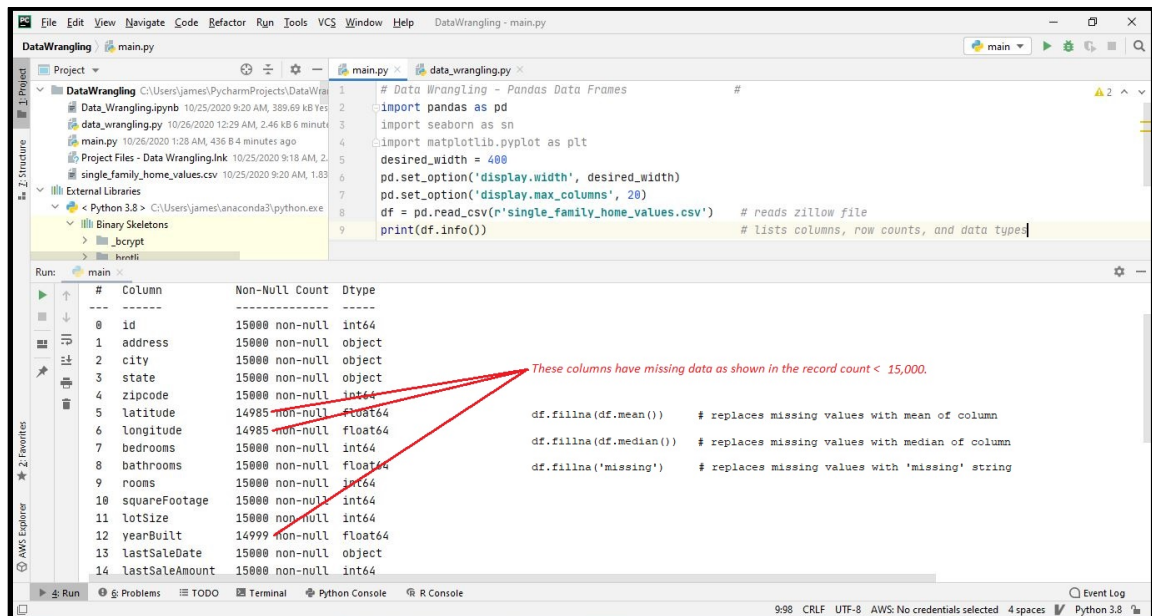
4) Once the csv file is loaded into a dataframe (df), we can begin data wrangling. Any of the following Python functions may be used with a df, once it has been assigned.

```
print(df.head(2))        # returns top two rows
print(df.tail(5))        # returns bottom five rows
print(type(df))          # returns type of df
print(df.shape)          # returns row & column count
```
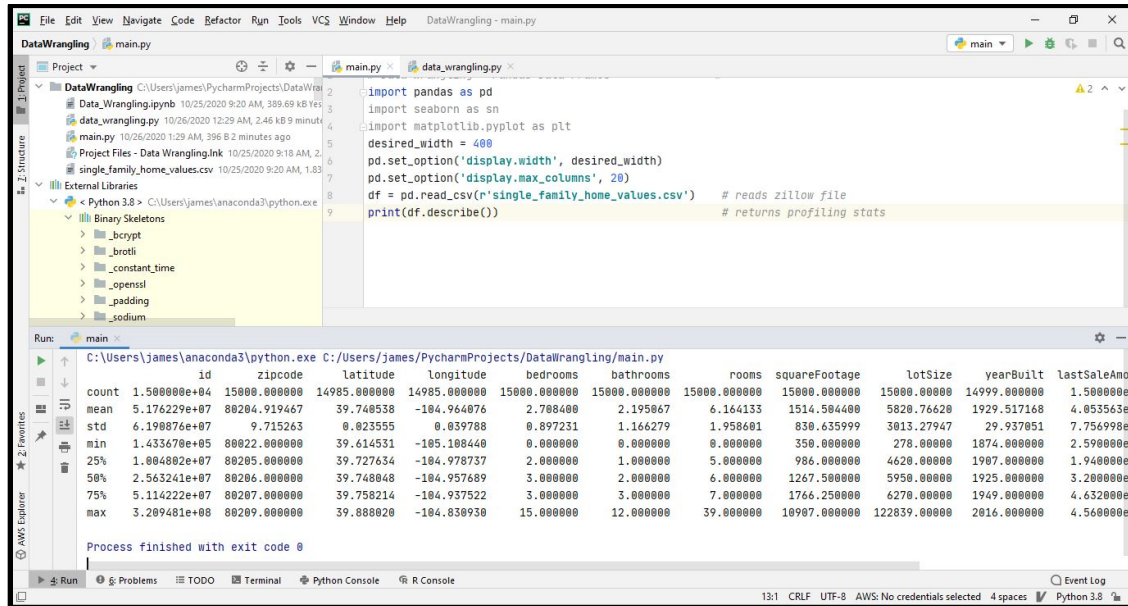


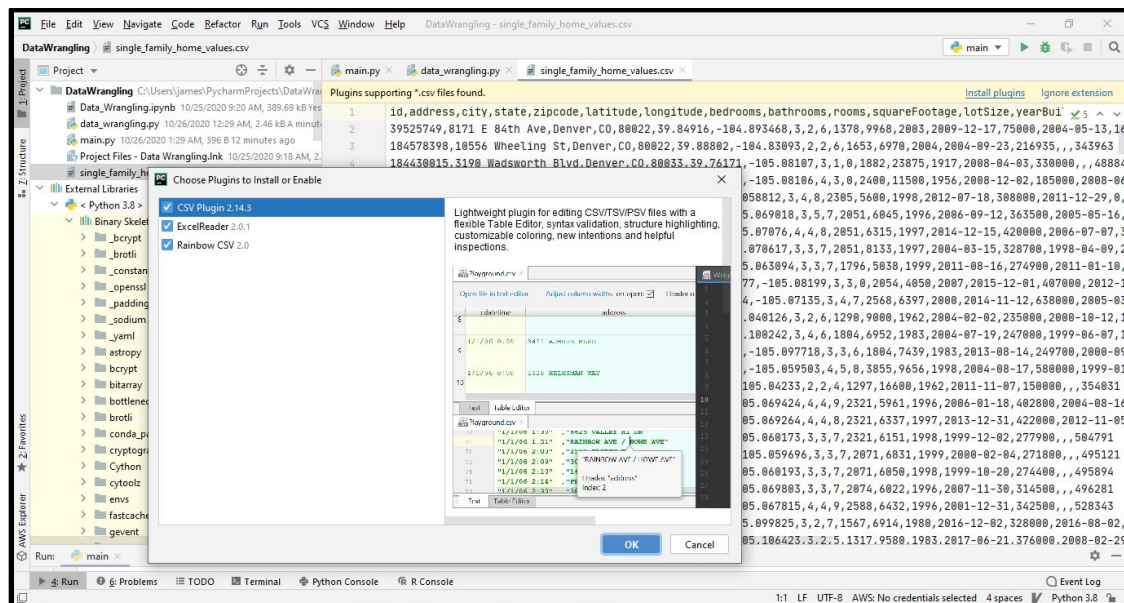*These images can be viewed at 400% magnification for details of the screen shots.*

```
print(df.info())        # returns column attributes
```



```
print(df.describe()) # returns profiling stats
```

5)  In PyCharm, print(df) will display the DataFrame.  You can also open the csv file in PyCharm by finding the file name ***single_family_home_values.csv*** in the project window.  Double click on the csv file name and it will open for you to few all 15,000 records.  You can also configure PyCharm to use a 3rd-party CSV viewer as shown in the example below.



6)  The pandas df functions df.shape, df.head(), df.tail(),  df.describe(), and df.shape may all be used as-is to return pandas dataframe details as illustrated in the previous screenshots. Additionally, these pandas df functions may also be used on individual column attributes.  The following examples illustrate theses pandas df functions for the attribute priorSaleAmount.

***df.priorSaleAmount.head(), df.priorSaleAmount.tail(), df.priorSaleAmount.shape***

*df.priorSaleAmount.describe(), df.priorSaleAmount.mean(), df.priorSaleAmount.median()*



7)  **Exercise deliverable:**  Construct three examples of pandas df functions in Python code.  Use either the single_family_home_values.csv or your data project csv.  (1)The first PyCharm screenshot should include Python code on importing pandas, reading the csv file, and displaying the attribute columns.  (2)The second PyCharm screenshot should be using the df.describe() function on an attribute or the entire dataframe.  (3)The third PyCharm screenshot  example should show the first two rows, last three rows, as well as the median of an attribute column. Capture the three PyCharm screenshots for use later in this tutorial.

## ✦ Dealing with missing data and outliers:

The pandas df function df.info() lists all the column attributes and the number of rows associated with each attribute.  As noted earlier, df.fillna() provides opportunity to compensate for missing data.  Outliers can also be identified in a bloxplot and then screened from being used in data analysis.  The boxplots of estimated_value data points created using the Python code below them show the impact of outliers.  Notice how the tail on the boxplot with outliers is troublesome:



**With Outliers**



**Without Outliers**

```
# Data Wrangling - Pandas Data Frames                    #
import pandas as pd
import numpy as np
import seaborn as sn
import matplotlib.pyplot as plt
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
df = pd.read_csv(r'single_family_home_values.csv')    # zillow file
print(df[['estimated_value', 'yearBuilt', 'priorSaleAmount']].tail(5))
print(df.estimated_value.shape, df.yearBuilt.shape, df.priorSaleAmount.shape)
sn.boxplot(df.estimated_value)
plt.show()
```

The following PyCharm screenshots creates a df2 dataframe to screen homes whose last sale was over $1,000,000 and/or estimated value is over $1,000,000.  The two example screenshots use an intersection of homes estimated below and (&)  previously sold below $1,000,000 as well as a union between homes estimated below $1,000,000 or (|)previously sold for less than $1,000,000.



*estimated_value and priorSalesAmount*



*estimated_value or priorSalesAmount*

1) **Exercise deliverable:** Construct an example screenshot of using pandas df functions in Python code to remove outliers in either the single_family_home_values.csv or your data project csv. If you are using the csv above, then choose another column attribute than priorSalesAmount. (1)Be sure to run a preliminary boxplot. (2) Remove the outliers in a new dataframe. (3) Run a new boxplot after the outliers removal. Save the PyCharm screenshot for use later in this tutorial.

## Simple Statistics in Python:

1) In addition to the pandas df function df.describe() reporting mean, standard deviation, minimum, maximum, percentiles (e.g., 25%, 50%, and 75%) metrics, the df.corr() function will provide simple correlation between all non-string column attributes. Review the following Python code:

```python
# Data Wrangling --- df2.corr()for one column vs. : (i.e., all columns    #
import pandas as pd
import numpy as np
import seaborn as sn
import matplotlib.pyplot as plt
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
df = pd.read_csv(r'single_family_home_values.csv')    # zillow file
df2 = df[(df.estimated_value<=1000000) | (df.lastSaleAmount<=1000000)]
print(df2.corr().loc['estimated_value', :].sort_values(ascending=False))

# Data Wrangling --- df2.corr()between all columns                        #
import matplotlib.pyplot as plt
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
df = pd.read_csv(r'single_family_home_values.csv')    # zillow file
df2 = df[(df.estimated_value<=1000000) | (df.lastSaleAmount<=1000000)]
print(df2.corr())
```

Note that in the first df2.corr() example above, the correlation is between estimated_value and all the other columns. The correlation results are sorted descending. The following screenshot is of the first df2.corr() correlation code example.



**Exercise deliverable:** Construct an example screenshot of using pandas df.corr() function in Python code to display the correlation between a column attribute and all other columns. Sort the correlation coefficients in decreasing order. Use either the single_family_home_values.csv or your data project csv. . Save the PyCharm screenshot for use later in this tutorial.

### ➕ Slicing a DataFrame:

Pandas df functions allow OLAP slicing.  Review the following Python code:
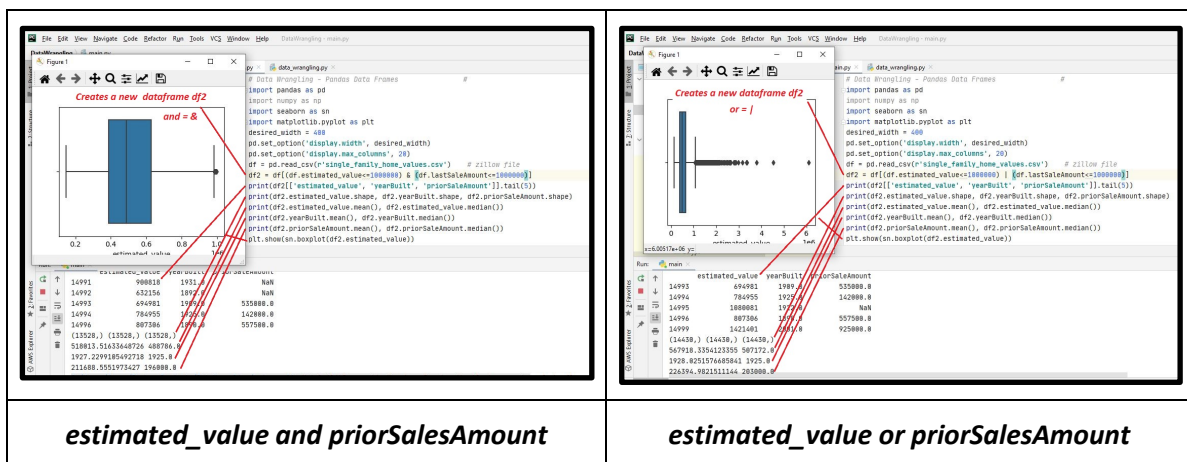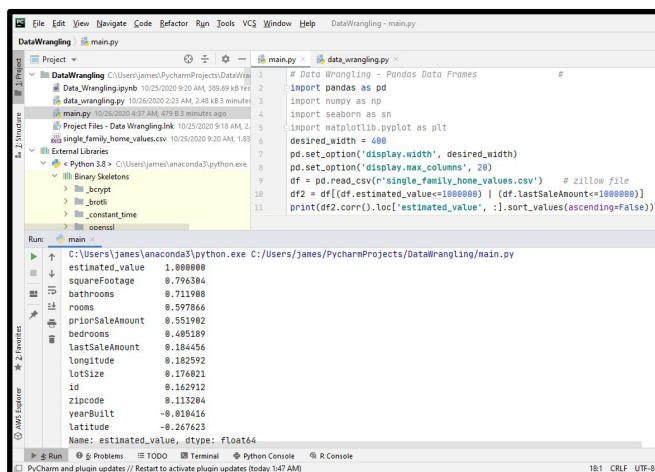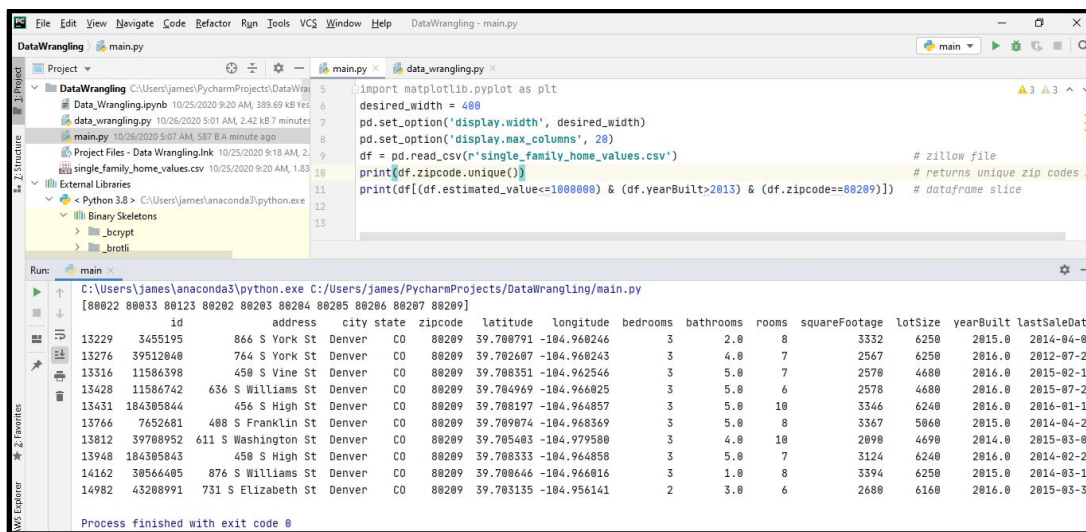
```python
# Data Wrangling - Pandas Data Frames                        #
import pandas as pd
import numpy as np
import seaborn as sn
import matplotlib.pyplot as plt
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
df = pd.read_csv(r'single_family_home_values.csv')
print(df.zipcode.unique())
print(df[(df.estimated_value<=1000000) & (df.yearBuilt>2013) & (df.zipcode==80209)])
```



**Exercise deliverable:** Construct an example screenshot of using pandas df functions in Python code to slice a dataframe by three or more column attributes.  Use either the single_family_home_values.csv or your data project csv.  If you choose the csv file above, do not use the exact same attributes.  The zip code may be reused.  Save the PyCharm screenshot for use later in this tutorial.

### ➕ Converting a Column, GroupBy(), and Merge():

Pandas df functions allow adding or dropping a column in the dataframe, grouping columns by column attributes, and merging dataframes.  The following series of Python code represents examples each:

1) ***Converting a Column:***  In PyCharm, the csv file column attribute lastSaleDate is a string object.  Converting the column attribute to a date format type has advantages in data mining and OLAP slices.  The following Python code and PyCharm screenshot illustrates how to add a converted column (i.e., lastSaleDate2) to a new df2.
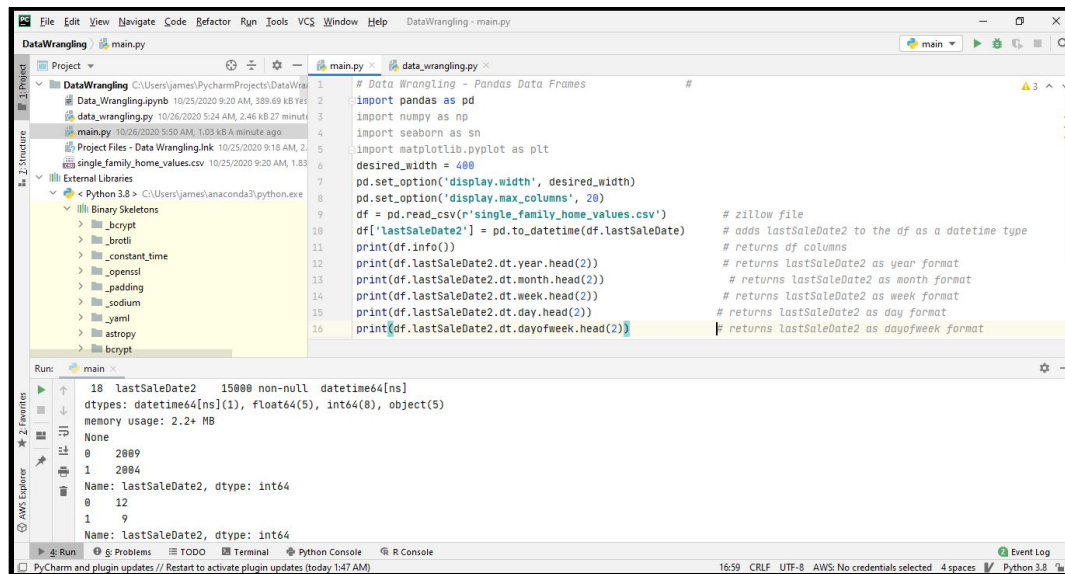
```python
# Data Wrangling - Pandas Data Frames                    #
import pandas as pd
import numpy as np
import seaborn as sn
import matplotlib.pyplot as plt
desired_width = 400
pd.set_option('display.width', desired_width)
```

```
pd.set_option('display.max_columns', 20)
df = pd.read_csv(r'single_family_home_values.csv')      # zillow file
df['lastSaleDate2'] = pd.to_datetime(df.lastSaleDate)  # adds to df as datetime type
print(df.info())                                        # returns df columns
print(df.lastSaleDate2.dt.year.head(2))        # returns lastSaleDate2 as year format
print(df.lastSaleDate2.dt.month.head(2))       # returns lastSaleDate2 as month format
print(df.lastSaleDate2.dt.week.head(2))        # returns lastSaleDate2 as week format
print(df.lastSaleDate2.dt.day.head(2))         # returns lastSaleDate2 as day format
print(df.lastSaleDate2.dt.dayofweek.head(2))   # returns as dayofweek format
```
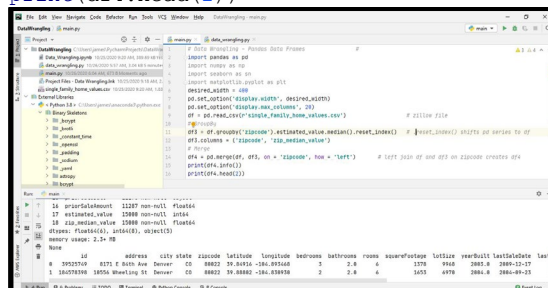


2) ***GroupBy & Merge:*** The pandas df function groupby() allows grouping and aggregate scalar values. In the cvs file, we can average the expected_value by zipcode into a new df3. The new df3 can then be merged with the old df to create a new df4. The Python code to perform this task and a screenshot of the code executed in PyCharm follows:

```
Data Wrangling – Pandas Data Frames                     #
import pandas as pd
import numpy as np
import seaborn as sn
import matplotlib.pyplot as plt
desired_width = 400
pd.set_option('display.width', desired_width)
pd.set_option('display.max_columns', 20)
df = pd.read_csv(r'single_family_home_values.csv')              # zillow file
# GroupBy
df3 = df.groupby('zipcode').estimated_value.median().reset_index() # .reset_index()
#shifts pd series to df
df3.columns = ('zipcode', 'zip_median_value')
# Merge
df4 = pd.merge(df, df3, on = 'zipcode', how = 'left')       # left join df and df3 on
#zipcode creates df4
print(df4.info())
print(df4.head(2))
```

3) ***Exercise deliverable:*** Construct an example screenshot of using pandas df functions in Python code for a df column conversion, groupby, and/or merge.  Use either the single_family_home_values.csv or your data project csv.  If you choose the csv file above, do not use the exact same attributes.  The zip code may be reused.  Save the PyCharm screenshot for use later in this tutorial.

### ✚ Deliverables from Exercises Above:

Take the noted deliverables (i.e., PyCharm screenshots) from each exercise above (i.e., total of 4 screenshots), insert the screenshots into a MS Word document as you label each screenshot deliverable from the corresponding exercise, save the document as one PDF, and submit the PDF to the Tutorial Seven Canvas Assignment uplink having the following qualities:

a. Your screenshot for each of the two exercises above is included.
b. Each screenshot is labeled appropriately for each exercise.
   i. Reading a CSV file into a Pandas DataFrame ………….. 3 screenshots
   ii. Dealing with missing data and outliers  …………………. 2 screenshots
   iii. Simple Statistics in Python  ……………………………………. 1 screenshot
   iv. Slicing a DataFrame  ……………………………………………. 1 screenshot
   v. Converting a Column, GroupBy(), and Merge() ………  1 screenshot
      Total PyCharm screenshots to submit ……………………  8 screenshots
c. All screenshots are legible output for each PyCharm exercise deliverable.