NOAA-CESSRST
CENTER FOR EARTH SYSTEM SCIENCES
AND REMOTE SENSING TECHNOLOGIES

NOAA
NATIONAL OCEANIC AND ATMOSPHERIC ADMINISTRATION
U.S. DEPARTMENT OF COMMERCE

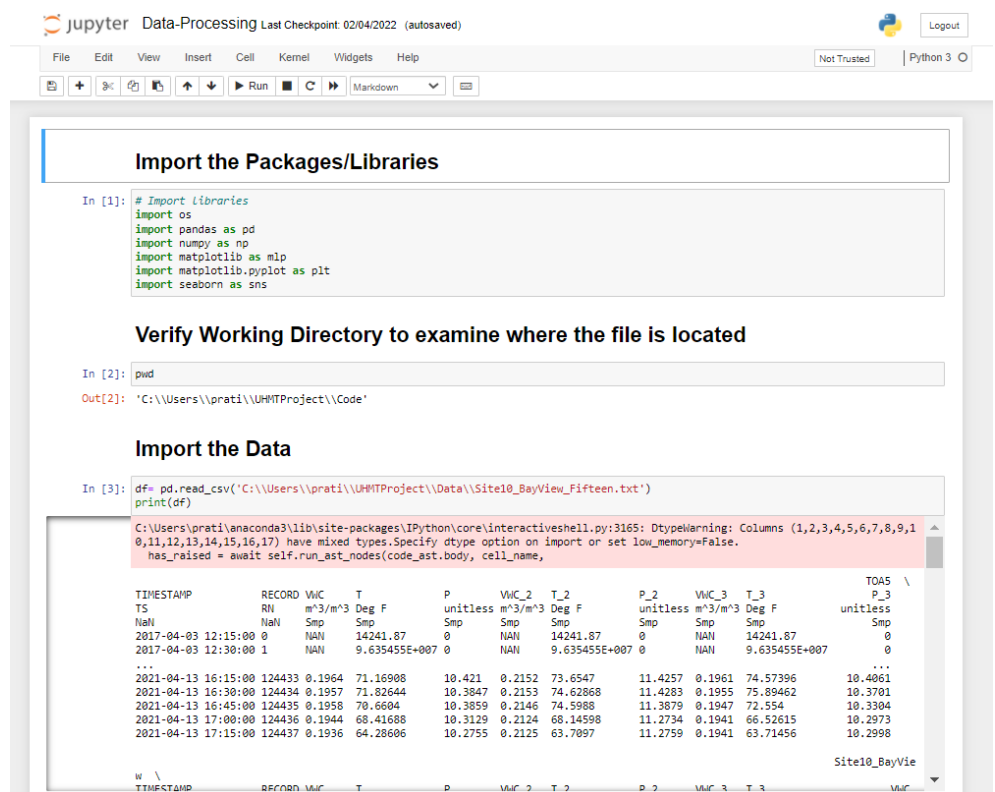Puja Roy                                                                                                    2/18/22

CET 4900 - OL60

## Internship Journal Entry #4

Throughout this week of my internship, I worked on analyzing large-scale datasets of real-time data collected from the New York Urban Hydro-Meteorological Testbed (NY-uHMT) weather station. Since there is 29+ large-scale datasets, it will take excessive time to analyze the data. The main objective of analyzing the large-scale datasets is to develop a Python script in Jupyter Notebook/Google Colaboratory to build a Weather App that automates the retrieval of weather stations and forecasts climate changes/natural hazards in NYC maps utilizing the New York Urban Hydro-Meteorological Testbed (NY-uHMT) data. Since data is not always organized and all over the place in Excel files/text files, it is necessary to clean, filter, examine and analyze the data to draw conclusions. My knowledge and skills of Python are being leveraged throughout my internship because I took EMT 1111 – Logic and Problem Solving during my freshman year of college. I gained exposure in learning how to program in Python and how it is mostly used for back-end programming. Most importantly, I am applying data science skills to my internship project because I am preparing the NY-uHMT data for analysis including cleansing, aggregating, and manipulating the data to visualize correlations and trends among various variables.

I downloaded all the data txt files that were collected by the NY-uHMT weather station. I used Python in Jupyter Notebook to import the data. For data analysis, I imported Python packages including pandas, numpy, matplotlib and seaborn. These packages are essential for reading and visualizing data files. Then, I printed the pwd command to verify the working directory in where the data files are located. Once the data was imported, I printed the dataframe which is declared df as a variable. The data for Bay View consists of 89,211 rows and 8 columns.



**Figure 1 – Python Script in Jupyter Notebook preparing for data analysis by importing Python packages and NY-uHMT data files**