# CS 487/519 Applied Machine Learning I
## Project 2: Open ML project
October 05, 2018

**Jithin Jacob Benjamin Jacob – 800681973**
**Joanna Augustine – 800656114**
**Group 9**

## MOTIVATION:

The key motivation behind this project is that carbon emissions contribute to climate change, which can have serious consequences for humans and their environment. According to the U.S. Environmental Protection Agency, carbon emissions, in the form of carbon dioxide, make up more than 80 percent of the greenhouse gases emitted in the United States. The burning of fossil fuels releases carbon dioxide and other greenhouse gases. These carbon emissions raise global temperatures by trapping solar energy in the atmosphere. So in this project we analyze a data on carbon footprint of individual's and using this data and machine learning techniques develop an algorithm that minimizes the carbon footprint of each individual while maintaining their quality of life.

## PROBLEM DEFINITION:

The problem of high carbon footprint in the environment is a major issue and has a great influence in the climatic changes in the world. So this problem is kept forth by Wells Fargo whose high priority is to promote environmental sustainability. As an initiative they have produced a data containing all the daily activities of individual customers. We are to analysis on the daily activities of individuals customers in a way to accelerate the transition to a low-carbon economy. This has to be achieved without compromising on their daily priorities and needs. They believe that taking individual actions can encourage the collective responsibility to achieve this goal. So using Machine Learning we are to develop a data product that would help in analyzing the data and help individuals to optimize the balance between their carbon footprint and the quality of life.

The ultimate goal would be to recommend an environment friendly change to the everyday actions without lessening the individuals' quality of life. The data gives a peak into the lives of 1,000 individuals who rated several everyday activities (taking a long shower, driving a car, etc.) on a scale of 1-100 based on how important those activities are to their daily lives. So at the end the data product should produce a computer data program to find quality substitutes for activities that are high carbon emitters without reducing the happiness and utility that the individuals in the data obtain from these activities.

## PROPOSED SOLUTION:

The solution that we are expecting from this analytical process is to refine the dataset in such a way that we can perform the basic three operations of loading the data, cleaning the data and thereby using the machine learning techniques to join the data which makes more relevance to the problem under discussion. So by the end of the problem we will get a more refined information on the dataset and thereby help in deciding on the alternatives that can be considered in order to arrive at our solution of low carbon footprint by the individuals.

**LINK TO DATASET:**

(Data set is also placed in our Github repository inside the stage3 folder)


**PROGRAM AND OUTPUTS:**

**(LOADING DATA)**

**DataLoader.py**

This jupyter notebook program is responsible to load the data from the excel file from both the sheets and just clean it up a little bit and save the data into csv files which is easier for loading in further analysis.

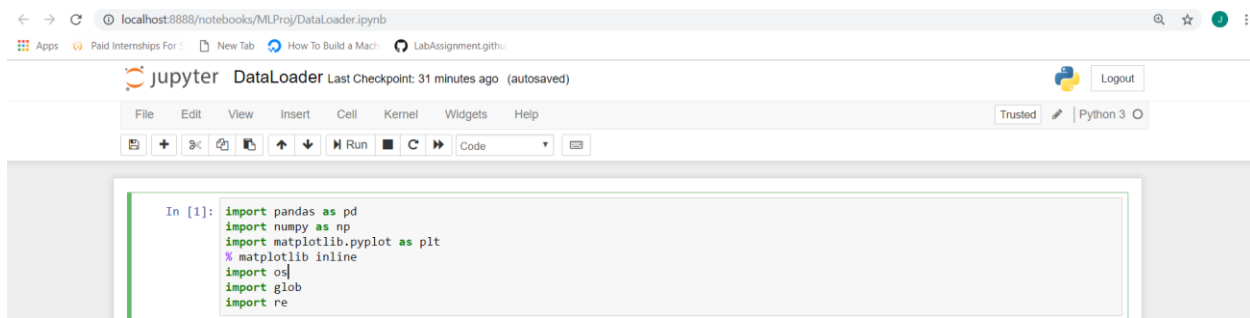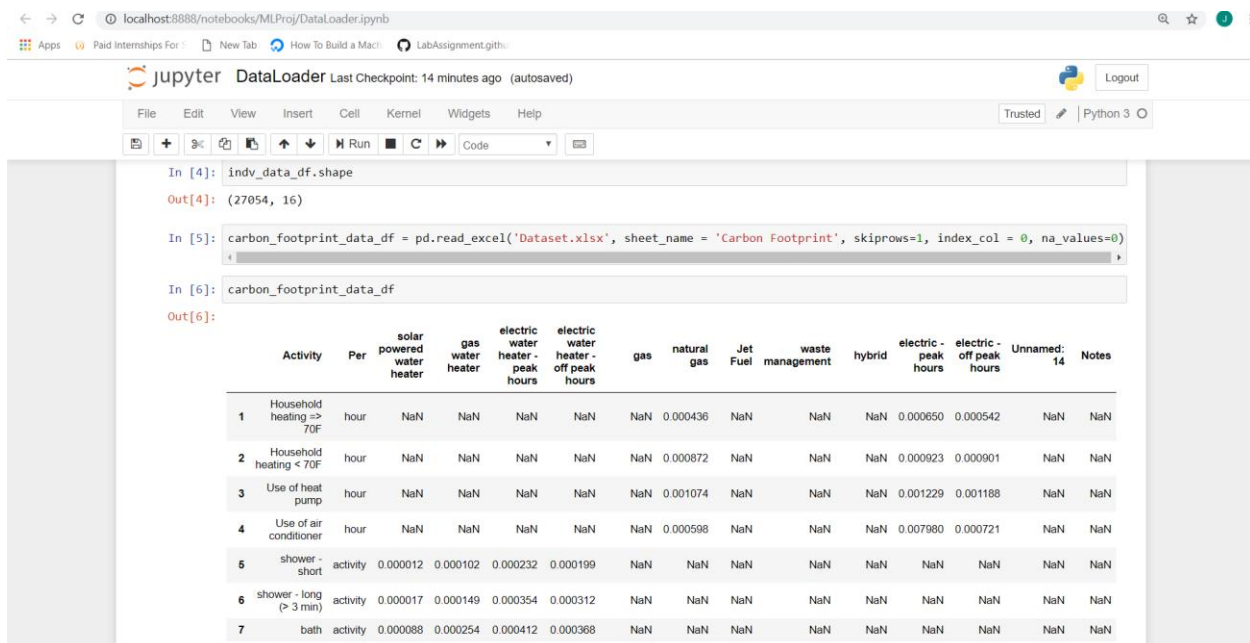**FIGURE 1: Importing the required libraries**



**FIGURE 2: Loading the Data from the Excel File**

```
In [7]: carbon_footprint_data_df.shape

Out[7]: (27, 15)
```



**FIGURE 3: Notebook After Performing Data Loader**