

# CS 487/519 Applied Machine Learning I

## Project 2: Open ML project – Stage 4

October 16, 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 to low carbon footprint of the individuals.

## LINK TO DATASET:

[https://www.mindsumo.com/contests/campus-analytics-challenge-2018?utm\\_campaign=send\\_drip\\_email&utm\\_source=mindsumo&utm\\_medium=email](https://www.mindsumo.com/contests/campus-analytics-challenge-2018?utm_campaign=send_drip_email&utm_source=mindsumo&utm_medium=email)

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

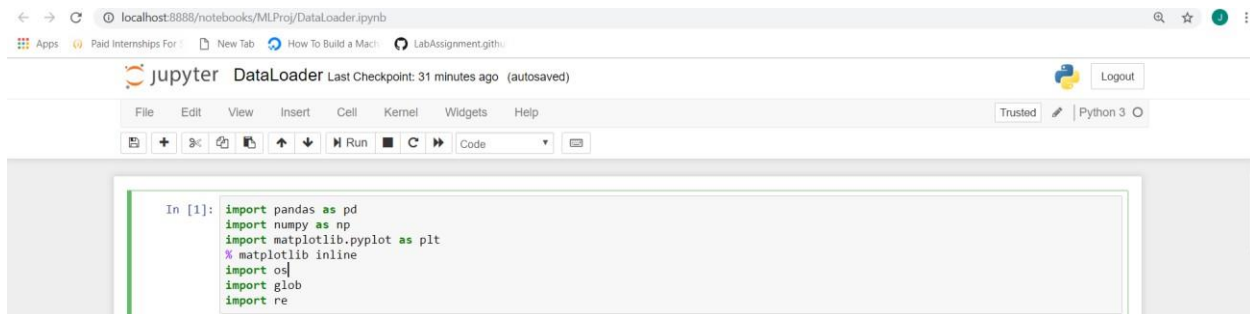
## PROGRAM AND OUTPUTS:

### (LOADING DATA)

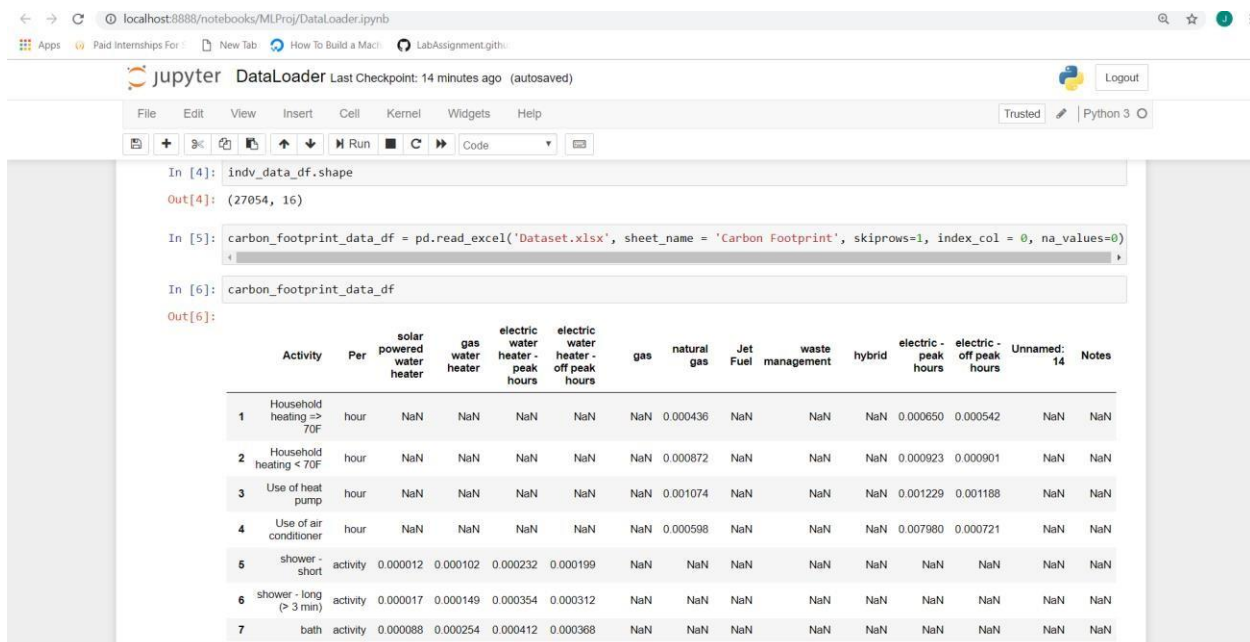
#### DataLoader.py

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)

In [8]: carbon_footprint_data_df_select_data = carbon_footprint_data_df.drop(['Unnamed: 14', 'Notes'], axis = 1)

In [9]: carbon_footprint_data_df_select_data.shape
Out[9]: (27, 13)

In [11]: carbon_footprint_data_df_select_data_table = pd.melt(carbon_footprint_data_df_select_data, id_vars=['Activity', 'Per'], value_var
C:\Users\d\Anaconda3\lib\site-packages\pandas\core\reshape\reshape.py:731: FutureWarning:
Passing list-likes to .loc or [] with any missing label will raise
KeyError in the future, you can use .reindex() as an alternative.

See the documentation here:
http://pandas.pydata.org/pandas-docs/stable/indexing.html#deprecate-loc-reindex-listlike
frame = frame.loc[:, id_vars + value_vars]
C:\Users\d\Anaconda3\lib\site-packages\pandas\core\indexing.py:1367: FutureWarning:
Passing list-likes to .loc or [] with any missing label will raise
KeyError in the future, you can use .reindex() as an alternative.

See the documentation here:
http://pandas.pydata.org/pandas-docs/stable/indexing.html#deprecate-loc-reindex-listlike
return self._getitem_tuple(key)

In [12]: carbon_footprint_data_df_select_data_table.shape

In [12]: carbon_footprint_data_df_select_data_table.shape
Out[12]: (351, 4)

In [13]: carbon_footprint_data_df_select_data_table_without_na = carbon_footprint_data_df_select_data_table.dropna(axis = 0)

In [14]: carbon_footprint_data_df_select_data_table_without_na.shape
Out[14]: (76, 4)

In [15]: carbon_footprint_data_df_select_data_table.to_csv('Carboon_Footprint_Table.csv')
carbon_footprint_data_df_select_data_table_without_na.to_csv('Carboon_Footprint_Table_out.csv')

In [ ]:
```

Figure 3: Notebook After Performing Data Loader

The screenshot shows the JupyterLab file browser interface. The breadcrumb path is 'localhost:8888/tree/MLProj'. The file list includes:

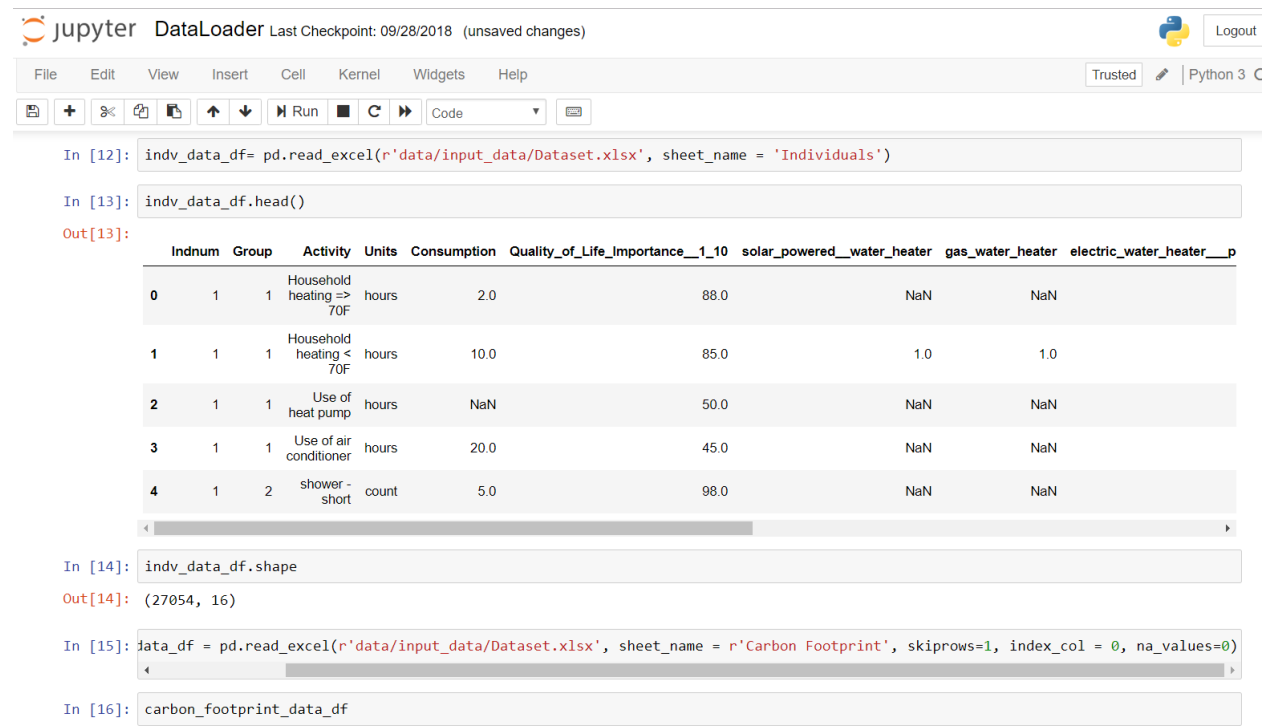
Name	Last Modified
..	seconds ago
Data Joiner.ipynb	Running 10 minutes ago
DataLoader.ipynb	Running 16 minutes ago
Carboon_Footprint_Table.csv	17 minutes ago
Carboon_Footprint_Table_out.csv	17 minutes ago
Dataset.xlsx	17 hours ago

## Stage 4 - Short Description:

In the stage 3 of the open MI project we had performed the Data Loading process where we had two csv files namely Carbon Footprint with na and Carbon Footprint without na which were used for the project for finding the low carbon footprint index and in this stage we put them into the input\_data and output\_data files respectively to categorize them for further analysis. **(Note: na is the missing data values denoted in pandas. It may also be abbreviated as Not Available)**

Secondly in this stage 4 we perform the Data Cleaning of the Individual Person's Data with respect to his daily activities which yield to higher carbon footprint in the environment.

**Figure 4: Loading the data from the original Dataset into the input\_data folder**



The screenshot shows a Jupyter Notebook titled 'DataLoader' with a last checkpoint of '09/28/2018 (unsaved changes)'. The interface includes a top bar with 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help' menus, along with a 'Trusted' status and 'Python 3' kernel selection. The notebook contains the following code and output:

```
In [12]: indv_data_df= pd.read_excel(r'data/input_data/Dataset.xlsx', sheet_name = 'Individuals')

In [13]: indv_data_df.head()

Out[13]:
```

	Indnum	Group	Activity	Units	Consumption	Quality_of_Life_Importance__1_10	solar_powered__water_heater	gas_water_heater	electric_water_heater__p
0	1	1	Household heating => 70F	hours	2.0	88.0	NaN	NaN	
1	1	1	Household heating < 70F	hours	10.0	85.0	1.0	1.0	
2	1	1	Use of heat pump	hours	NaN	50.0	NaN	NaN	
3	1	1	Use of air conditioner	hours	20.0	45.0	NaN	NaN	
4	1	2	shower - short	count	5.0	98.0	NaN	NaN	

```
In [14]: indv_data_df.shape

Out[14]: (27054, 16)

In [15]: data_df = pd.read_excel(r'data/input_data/Dataset.xlsx', sheet_name = r'Carbon Footprint', skiprows=1, index_col = 0, na_values=0)

In [16]: carbon_footprint_data_df
```

**Figure 5: Extracting data and giving the sheet a name such as Carbon Footprint and displaying the data.**

Out[16]:

	Activity	Per	solar powered water heater	gas water heater	electric water heater - peak hours	electric water heater - off peak hours	gas	natural gas	Jet Fuel	waste management	hybrid	electric - peak hours	electric - off peak hours	Unnamed: 14	Notes
1	Household heating => 70F	hour	NaN	NaN	NaN	NaN	NaN	0.000436	NaN	NaN	NaN	0.000650	0.000542	NaN	NaN
2	Household heating < 70F	hour	NaN	NaN	NaN	NaN	NaN	0.000872	NaN	NaN	NaN	0.000923	0.000901	NaN	NaN
3	Use of heat pump	hour	NaN	NaN	NaN	NaN	NaN	0.001074	NaN	NaN	NaN	0.001229	0.001188	NaN	NaN
4	Use of air conditioner	hour	NaN	NaN	NaN	NaN	NaN	0.000598	NaN	NaN	NaN	0.007980	0.000721	NaN	NaN
5	shower - short	activity	0.000012	0.000102	0.000232	0.000199	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
6	shower - long (> 3 min)	activity	0.000017	0.000149	0.000354	0.000312	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
7	bath	activity	0.000088	0.000254	0.000412	0.000368	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
8	wash-up	activity	0.000004	0.000042	0.000067	0.000055	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
9	use of dishwasher	activity	0.000025	0.000165	0.000398	0.000311	NaN	NaN	NaN	NaN	NaN	0.000084	0.000078	NaN	NaN
10	use of clothes washer	activity	0.000033	0.000199	0.000433	0.000382	NaN	0.000154	NaN	NaN	NaN	0.000102	0.000093	NaN	NaN
11	use of clothes dryer	activity	NaN	NaN	NaN	NaN	NaN	0.000187	NaN	NaN	NaN	0.000132	0.000122	NaN	NaN
	use of														

**Figure 6: Dropping the columns which contains the data which are less or not relevant to carbon footprint impact and defining carbon\_footprint\_data\_df\_select\_data\_ as the refined value, here the data with na is considered.**

This method is more of a refining process to eliminate the unwanted columns.  
( Note: The warning here is to use .reindex() as an alternative.)

```

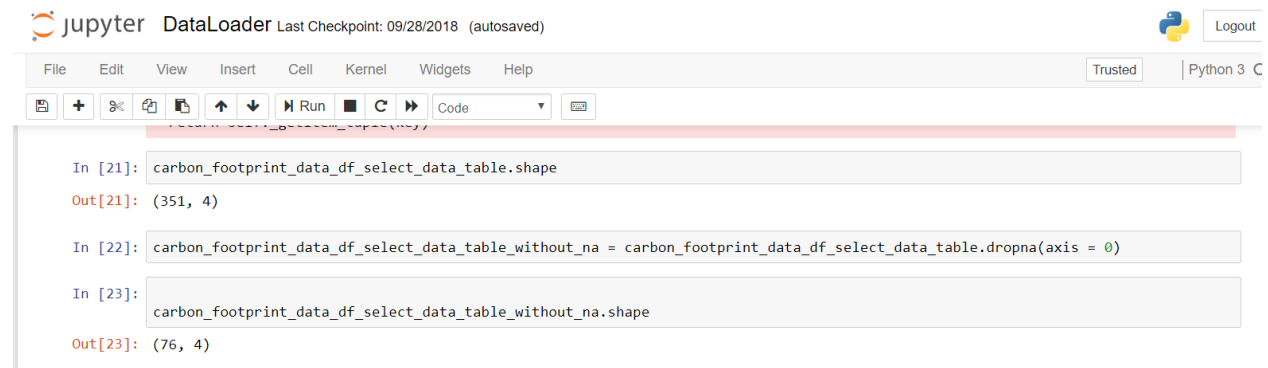
jupyter DataLoader Last Checkpoint: 09/28/2018 (autosaved)
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 C
In [17]: carbon_footprint_data_df.shape
Out[17]: (27, 15)
In [18]: carbon_footprint_data_df_select_data = carbon_footprint_data_df.drop(['Unnamed: 14', 'Notes'], axis = 1)
In [19]: carbon_footprint_data_df_select_data.shape
Out[19]: (27, 13)
In [20]: c(carbon_footprint_data_df_select_data, id_vars=['Activity', 'Per'], value_vars=list(carbon_footprint_data_df.columns.values)[2:])
C:\Users\d\Anaconda3\lib\site-packages\pandas\core\reshape\reshape.py:731: FutureWarning:
Passing list-likes to .loc or [] with any missing label will raise
KeyError in the future, you can use .reindex() as an alternative.

See the documentation here:
http://pandas.pydata.org/pandas-docs/stable/indexing.html#deprecate-loc-reindex-listlike
frame = frame.loc[:, id_vars + value_vars]
C:\Users\d\Anaconda3\lib\site-packages\pandas\core\indexing.py:1367: FutureWarning:
Passing list-likes to .loc or [] with any missing label will raise
KeyError in the future, you can use .reindex() as an alternative.

See the documentation here:
http://pandas.pydata.org/pandas-docs/stable/indexing.html#deprecate-loc-reindex-listlike
return self._getitem_tuple(key)

```

**Figure 7: Dropping the columns which are not required and defining carbon\_footprint\_data\_df\_select\_data\_table\_without\_na as columns where axis=0 and where na is not present**



```

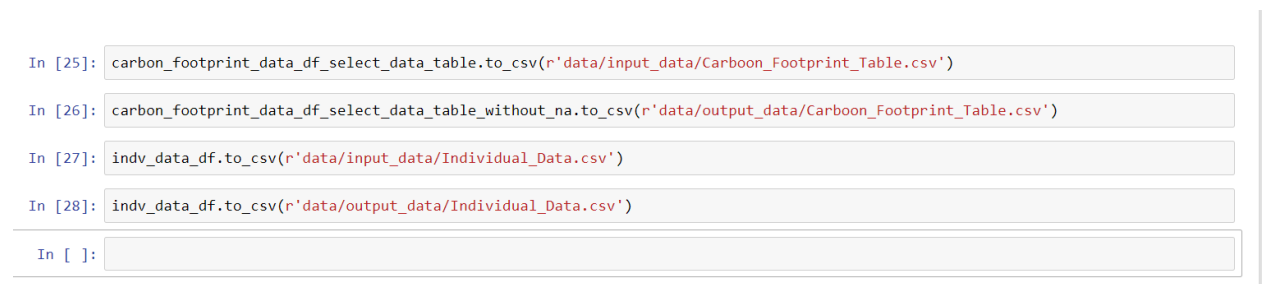
In [21]: carbon_footprint_data_df_select_data_table.shape
Out[21]: (351, 4)

In [22]: carbon_footprint_data_df_select_data_table_without_na = carbon_footprint_data_df_select_data_table.dropna(axis = 0)

In [23]: carbon_footprint_data_df_select_data_table_without_na.shape
Out[23]: (76, 4)

```

**Figure 8: Exporting the data into csv files and saving them in the input\_data and output\_data folders**



```

In [25]: carbon_footprint_data_df_select_data_table.to_csv(r'data/input_data/Carboon_Footprint_Table.csv')

In [26]: carbon_footprint_data_df_select_data_table_without_na.to_csv(r'data/output_data/Carboon_Footprint_Table.csv')

In [27]: indiv_data_df.to_csv(r'data/input_data/Individual_Data.csv')

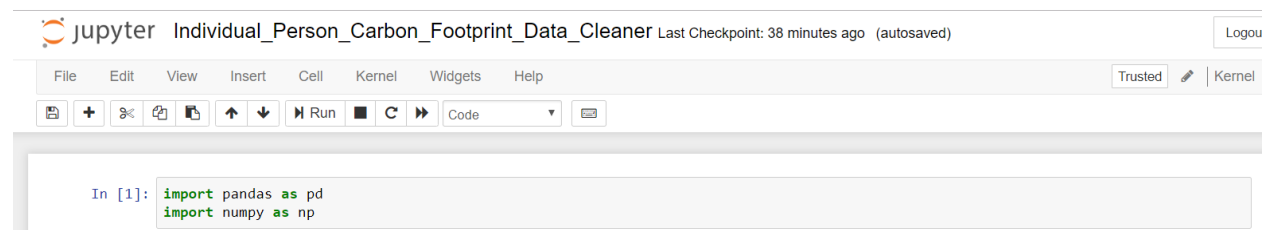
In [28]: indiv_data_df.to_csv(r'data/output_data/Individual_Data.csv')

In [ ]:

```

## Individual Person Carbon Footprint Data Cleaner

**Figure 9: First we import the required libraries**



```

In [1]: import pandas as pd
import numpy as np

```

Secondly we load the data from our first and main Dataset which is “Dataset.xlsx” and we output the values as shown in the figure below.

Jupyter Individual\_Person\_Carbon\_Footprint\_Data\_Cleaner Last Checkpoint: 41 minutes ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Kernel C

In [2]: individuals\_carbon\_footprint\_df = pd.read\_excel(r'data/input\_data/Dataset.xlsx', sheet\_name='Individuals')

In [3]: individuals\_carbon\_footprint\_df.head()

Out[3]:

	Indnum	Group	Activity	Units	Consumption	Quality_of_Life_Importance__1_10	solar_powered__water_heater	gas_water_heater	electric_water_heater__p
0	1	1	Household heating => 70F	hours	2.0	88.0	NaN	NaN	
1	1	1	Household heating < 70F	hours	10.0	85.0	1.0	1.0	
2	1	1	Use of heat pump	hours	NaN	50.0	NaN	NaN	
3	1	1	Use of air conditioner	hours	20.0	45.0	NaN	NaN	
4	1	2	shower - short	count	5.0	98.0	NaN	NaN	

In [4]: individuals\_carbon\_footprint\_df.shape

Out[4]: (27054, 16)

**Figure 10: Augmenting the data**

Thirdly we need to Augment the data which is basically converting the data from a 2 Dimensional Table to a 1 Dimensional Table by comparing the data of every individual vs the type of resource that he/she is using.

Jupyter Individual\_Person\_Carbon\_Footprint\_Data\_Cleaner Last Checkpoint: 43 minutes ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Kernel C

In [5]: individuals\_carbon\_footprint\_df.columns.values[6:], var\_name='Name of Resource Used', value\_name='Amount of Resource Used per Unit')

In [6]: individuals\_carbon\_footprint\_table.head()

Out[6]:

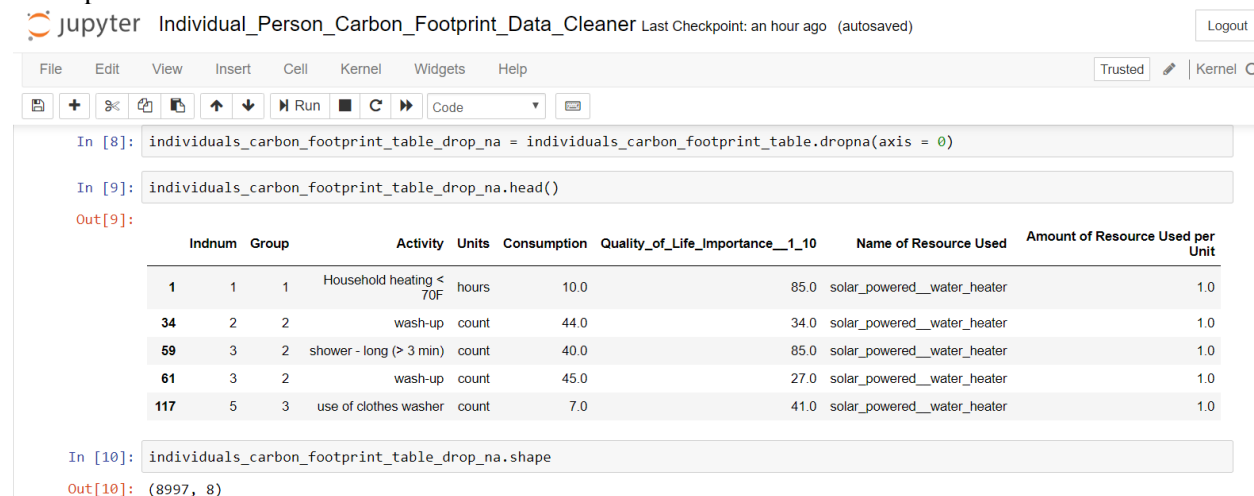
	Indnum	Group	Activity	Units	Consumption	Quality_of_Life_Importance__1_10	Name of Resource Used	Amount of Resource Used per Unit
0	1	1	Household heating => 70F	hours	2.0	88.0	solar_powered__water_heater	NaN
1	1	1	Household heating < 70F	hours	10.0	85.0	solar_powered__water_heater	1.0
2	1	1	Use of heat pump	hours	NaN	50.0	solar_powered__water_heater	NaN
3	1	1	Use of air conditioner	hours	20.0	45.0	solar_powered__water_heater	NaN
4	1	2	shower - short	count	5.0	98.0	solar_powered__water_heater	NaN

In [7]: individuals\_carbon\_footprint\_table.shape

Out[7]: (270540, 8)

**Figure 11: Removing the NaN values and replacing it with 0.0**

The fourth part deals with all the NaN values which were outputted after the augmenting of the data has taken place.



Jupyter Notebook interface for 'Individual\_Person\_Carbon\_Footprint\_Data\_Cleaner'. The notebook shows two code cells and their outputs.

```
In [8]: individuals_carbon_footprint_table_drop_na = individuals_carbon_footprint_table.dropna(axis = 0)
```

```
In [9]: individuals_carbon_footprint_table_drop_na.head()
```

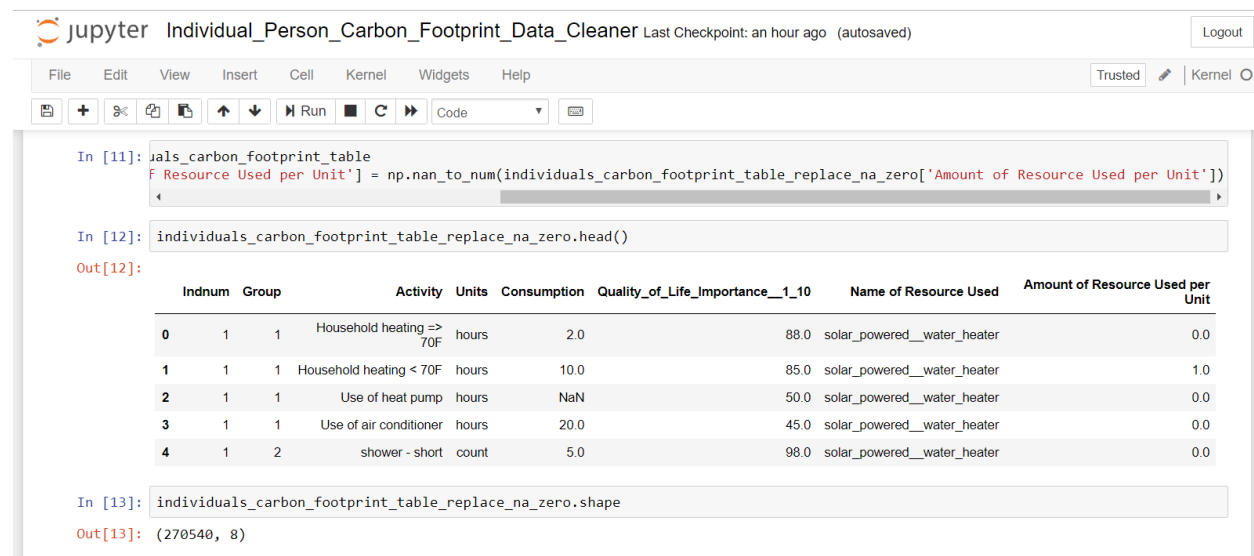
Out[9]:

	Indnum	Group	Activity	Units	Consumption	Quality_of_Life_Importance__1_10	Name of Resource Used	Amount of Resource Used per Unit
	1	1	Household heating < 70F	hours	10.0	85.0	solar_powered__water_heater	1.0
34	2	2	wash-up	count	44.0	34.0	solar_powered__water_heater	1.0
59	3	2	shower - long (> 3 min)	count	40.0	85.0	solar_powered__water_heater	1.0
61	3	2	wash-up	count	45.0	27.0	solar_powered__water_heater	1.0
117	5	3	use of clothes washer	count	7.0	41.0	solar_powered__water_heater	1.0

```
In [10]: individuals_carbon_footprint_table_drop_na.shape
```

Out[10]: (8997, 8)

Next we replace the value of NaN with 0.0 wherever NaN is found



Jupyter Notebook interface for 'Individual\_Person\_Carbon\_Footprint\_Data\_Cleaner'. The notebook shows two code cells and their outputs.

```
In [11]: individuals_carbon_footprint_table_replace_na_zero['Amount of Resource Used per Unit'] = np.nan_to_num(individuals_carbon_footprint_table_replace_na_zero['Amount of Resource Used per Unit'])
```

```
In [12]: individuals_carbon_footprint_table_replace_na_zero.head()
```

Out[12]:

	Indnum	Group	Activity	Units	Consumption	Quality_of_Life_Importance__1_10	Name of Resource Used	Amount of Resource Used per Unit
0	1	1	Household heating => 70F	hours	2.0	88.0	solar_powered__water_heater	0.0
1	1	1	Household heating < 70F	hours	10.0	85.0	solar_powered__water_heater	1.0
2	1	1	Use of heat pump	hours	NaN	50.0	solar_powered__water_heater	0.0
3	1	1	Use of air conditioner	hours	20.0	45.0	solar_powered__water_heater	0.0
4	1	2	shower - short	count	5.0	98.0	solar_powered__water_heater	0.0

```
In [13]: individuals_carbon_footprint_table_replace_na_zero.shape
```

Out[13]: (270540, 8)

Finally we save the tables to the input\_data and output\_data locations

**Figure 12: Saving the data to input\_data and output\_data locations after dropping na and replacing NaN by 0.0**

```
In [14]: individuals_carbon_footprint_table_drop_na.to_csv(r'data/output_data/Individuals_Carbon_Footprint_NA_Dropped.csv', index=False)
         individuals_carbon_footprint_table_replace_na_zero.to_csv(r'data/output_data/Individuals_Carbon_Footprint_NA_Zeroed.csv', index=False)
```

Our notebook finally looks like this:



jupyter

Logout

FilesRunningClusters

Select items to perform actions on them.

UploadNew↺

0 / MLProj / data / input\_data

Name↑Last Modified

..seconds ago

Individual\_Data.csvan hour ago

Dataset.xlsx18 days ago

Carboon\_Footprint\_Table.csvan hour ago

jupyter

Logout

FilesRunningClusters

Select items to perform actions on them.

UploadNew↺

0 / MLProj / data / output\_data

Name↑Last Modified

..seconds ago

Individuals\_Carbon\_Footprint\_NA\_Zeroed.csv15 minutes ago

Individuals\_Carbon\_Footprint\_NA\_Dropped.csv15 minutes ago

Individual\_Data.csvan hour ago

Carboon\_Footprint\_Table.csvan hour ago

## RESULT:

The above analysis have been made and the results have been shown in the figures above. These are to be updated in the future.