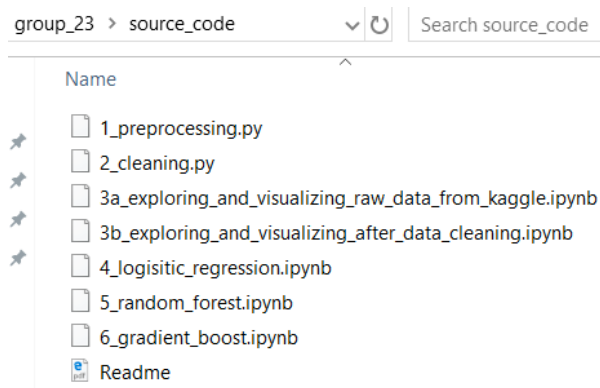


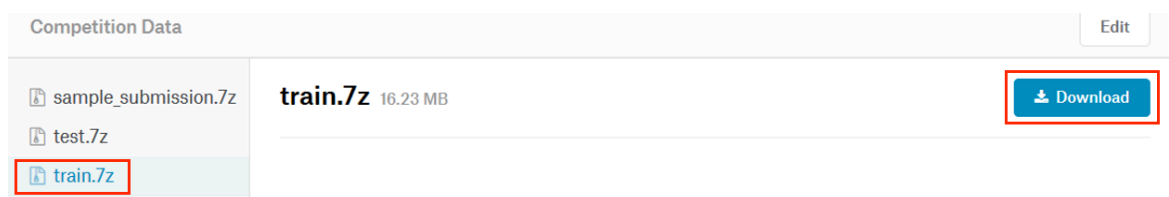
## Source code structure:

All the code files needed to run the model is shown in the screenshot below:



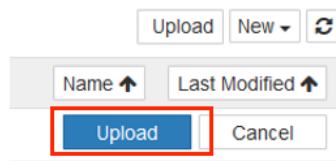
## Steps to run the code:

1. Download data file “**train.csv**” from <https://www.kaggle.com/c/porto-seguro-safe-driver-prediction/data>. Select train.7z and then click Download button. Save the file in a directory of your choice, say “c:\ids”



2. Unzip this file and extract here to get a file named “**train.csv**”. The path for our case should look like “c:\ids\train.csv”
3. Launch your Python IDE and then open file **1\_preprocessing.py**.
  - a. Search for “train = pd.read\_csv” and replace the directory where your **train.csv** file exists (“c:\ids\train.csv” in our case)
  - b. Search for “rows\_0\_1.to\_csv” and replace the directory where you want “**train\_u.csv**” file to exist (say “c:\ids\train\_u.csv” in our case)
4. Run the file. After the code is run, you will find the file here: “c:\ids\train\_u.csv”. Close file **1\_preprocessing.py**.
5. Open **2\_cleaning.py**.
  - a. Search for “train = pd.read\_csv” and replace the directory where your “**train\_u.csv**” file exists (“c:\ids\train\_u.csv” in our case)
  - b. Search for “train\_final\_df.to\_csv” and replace the directory where you want “**train\_final.csv**” file to exist (say “c:\ids\train\_final.csv” in our case)
6. Run the file. After the code is run, you will find the file here: “c:\ids\train\_final.csv”. Close file **2\_cleaning.py**.

- Open your Jupyter Notebook. Click "Upload" on the top-right corner. Select file **"3a\_exploring\_and\_visualizing\_raw\_data\_from\_kaggle.ipynb"** and open. Click Upload button again.



You will now find this file listed in your directory. Click on this file. Search for "train = pd.read\_csv" and replace the directory where your **train.csv** file exists ("c:\ids\train.csv" in our case). Click on "Cell" on the menu bar and then click "Run All". This will run the entire code set. You can browse through the outputs. This file details the feature drop analysis and missing data imputation.

- For data visualization and analysis, upload file **"3b\_exploring\_and\_visualizing\_after\_data\_cleaning.ipynb"** in Jupyter Notebook as per #7 above. Click on this file. Search for "train = pd.read\_csv" and replace the directory where your **train\_final.csv** file exists ("c:\ids\train\_final.csv" in our case). Click on "Cell" on the menu bar and then click "Run All". This will run the entire code set. You can browse through the outputs. This file details visual data analysis.
- For running Logistic Regression, upload file **"4\_logistic\_regression.ipynb"** in Jupyter Notebook as per #7 above. Click on this file. Search for "train = pd.read\_csv" and replace the directory where your **train\_final.csv** file exists ("c:\ids\train\_final.csv" in our case). Click on "Cell" on the menu bar and then click "Run All". This will run the entire code set. You can browse through the outputs.
- For running Random Forest Classifier, upload file **"5\_random\_forest.ipynb"** in Jupyter Notebook as per #7 above. Click on this file. Search for "train = pd.read\_csv" and replace the directory where your **train\_final.csv** file exists ("c:\ids\train\_final.csv" in our case). Click on "Cell" on the menu bar and then click "Run All". This will run the entire code set. You can browse through the outputs.
- For running Gradient Boost Classifier, upload file **"6\_gradient\_boost.ipynb"** in Jupyter Notebook as per #7 above. Click on this file. Search for "train = pd.read\_csv" and replace the directory where your **train\_final.csv** file exists ("c:\ids\train\_final.csv" in our case). Click on "Cell" on the menu bar and then click "Run All". This will run the entire code set. You can browse through the outputs.