# USER MANUAL GUIDE FOR CREATING A LOGISTIC MODEL

# Instructions to running GLM GUI
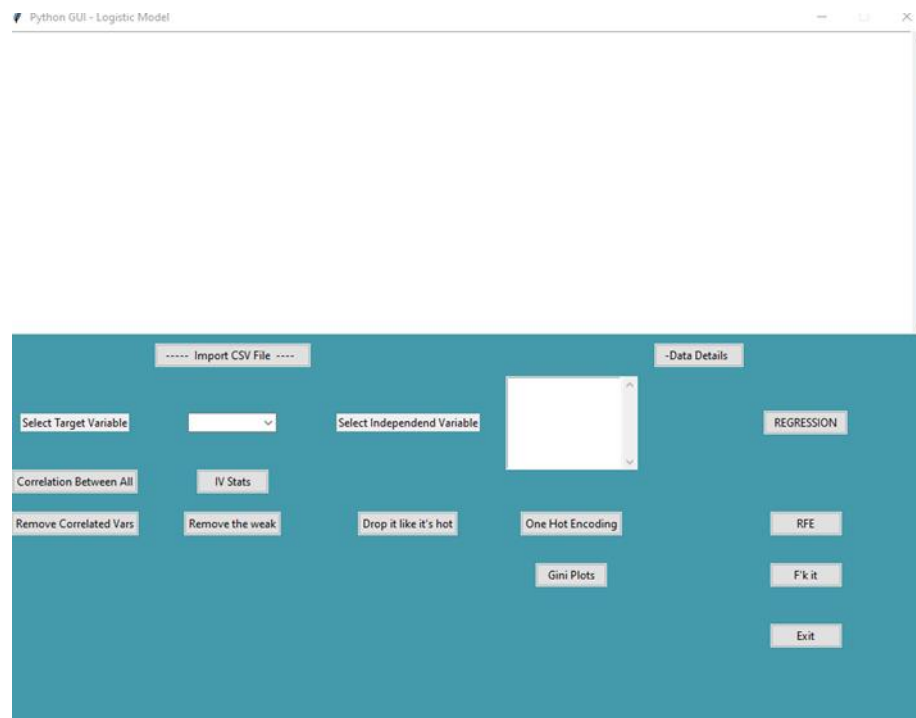
1) First run all of the libraries

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Tue Jul  9/7/2020
4
5 @author: Alun Brain (Dr. Brain Stats)
6 """
7
8 import tkinter as tk
9 from tkinter import ttk
10 from tkinter import *
11 from tkinter import scrolledtext
12 from tkinter import Tk
13 from tkinter.filedialog import askopenfilename
14 import seaborn as sns
15 import statsmodels.api as sm
16 #import statsmodels.formula.api as sm
17
18 import pandas as pd
19 import numpy as np
20 from scipy import stats
21 import matplotlib.pyplot as plt
22 from tkinter import filedialog
23
24
25 import pandas.core.algorithms as algos
26 from pandas import Series
27 import re
28 import traceback
29
30 import matplotlib
31 matplotlib.use("TkAgg")
32 from matplotlib.backends.backend_tkagg import ( FigureCanvasTkAgg, NavigationToolbar2Tk)
33 from matplotlib.figure import Figure
34
35
36 from sklearn.metrics import roc_auc_score
37 from sklearn.metrics import roc_curve
38 from sklearn.model_selection import train_test_split
39 #reduction using RFE
40 from sklearn.feature_selection import RFE
41 from sklearn.linear_model import LogisticRegression
42
```
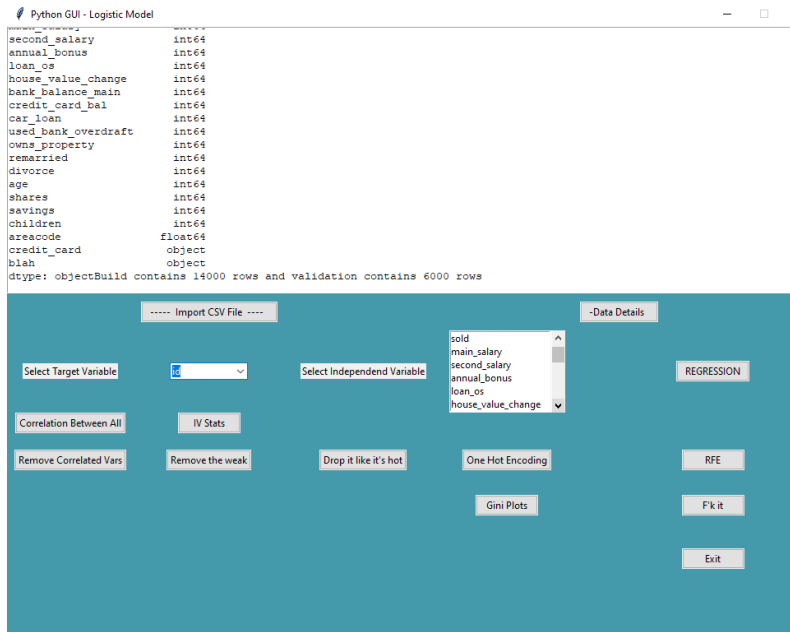
2) Run all the code between:
   a. Tk().withdraw()
   b. root.mainloop()
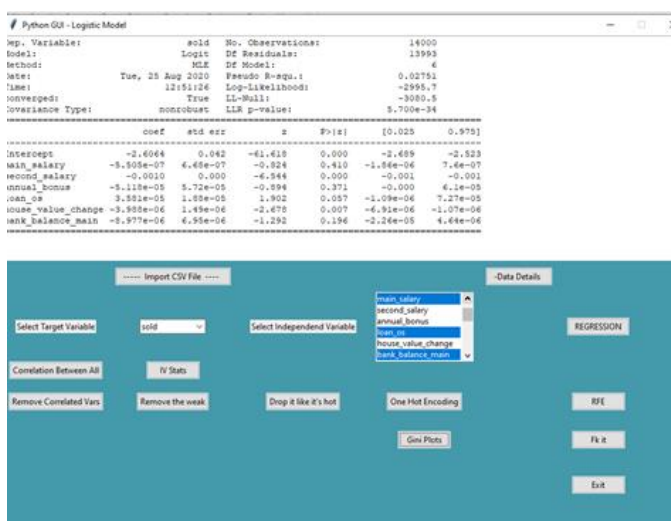3) You get the GUI

# Using the GUI – Logistic regression

1) Import CSV file – click this to get your CSV file



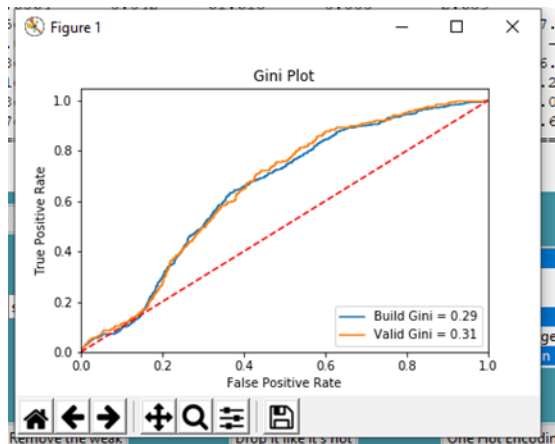This will tell you information about your data and split you data in70/30 split, ashown

2) Create Initial Multi-linear regression
   - Select Target Variable
     - This can easily be changed, we will be using **sold**
   - Select Independent Variable
     - Click on each variable you wish to model, in this scenario , main_salary, loan_os andbank_balance_main
   - Then press **REGRESSION**

Your model is in the screen.

Clicking on Gini-plot gives you



It applies the model on to your holdout, so a buildb holdout comparison can be made.

## Extra Funcationality

## Correlation



Will provide all of your correlated statistics.

Will also remove any correlated values based on 0.7 cut-off

## IV stats



Will provide all of your IV statistics.

Will also remove any variables with an IV lower than 0.1

## Drop it like its hot



Removes those categorical variables with over 10 bins (easily amendable). This is recommended before you use **One Hot encoding**

## One Hot Encoding



Create dummy variables on categorical variables, but removes one bin from each variable so it can not suffer from collinearity.

## RFE



Provides you with your top performing variables, but now variable reduction is done as this part is subjective.

# JUST GIVE ME A MODEL BUTTON



This conducts variable reduction, one hot encoding, model build and graph in one go



And so on…



And on….

Python GUI - Logistic Model

ime:                    13:08:53   Log-Likelihood:              -2494.4
inverged:                    True   LL-Null:                     -3080.5
ivariance Type:         nonrobust   LLR p-value:             1.318e-246
================================================================================
                    coef    std err         z      P>|z|     [0.025     0.975]
--------------------------------------------------------------------------------
itercept         -3.7381      0.074   -50.374      0.000     -3.883     -3.593
innual_bonus    -4.592e-05   4.92e-05    -0.933      0.351     -0.000   5.06e-05
ink_balance_main  7.969e-07   1.16e-06     0.689      0.491  -1.47e-06   3.07e-06
ildren            0.3462      0.022    16.063      0.000      0.304      0.388
.vorce            1.6659      0.108    15.394      0.000      1.454      1.878
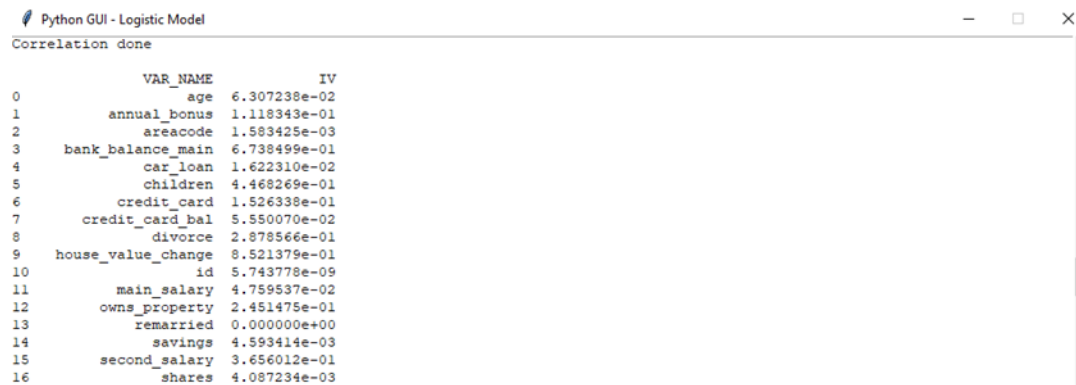juse_value_change -1.385e-06  1.22e-06    -1.136      0.256  -3.77e-06   1.01e-06
ins_property     -1.7440      0.165   -10.553      0.000     -2.068     -1.420
cond_salary      -0.0007      0.000    -5.010      0.000     -0.001     -0.000
ied_bank_overdraft  1.2996      0.096    13.518      0.000      1.111      1.488
edit_card_Y       0.4100      0.097     4.228      0.000      0.220      0.600
================================================================================