

Import all necessary open source libraries

```
In [1]: import pandas as pd
import numpy as np
from sklearn import preprocessing
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
%matplotlib inline
plt.rcParams['figure.figsize'] = (10, 5)
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report
from sklearn.metrics import roc_curve, roc_auc_score
from sklearn.metrics import confusion_matrix
import seaborn as sns
sns.set(style='white')
sns.set(style='whitegrid', color_codes=True)
import statsmodels.api as sm

In [2]: pip install -U imbalanced-learn

Requirement already satisfied: imbalanced-learn in /Users/divyabastola/opt/anaconda3/lib/python3.8/site-packages (0.10.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in /Users/divyabastola/opt/anaconda3/lib/python3.8/site-packages (from imbalanced-learn) (2.1.0)
Requirement already satisfied: joblib>=1.1.1 in /Users/divyabastola/opt/anaconda3/lib/python3.8/site-packages (from imbalanced-learn) (1.2.0)
Requirement already satisfied: scikit-learn>=0.22.0 in /Users/divyabastola/opt/anaconda3/lib/python3.8/site-packages (from imbalanced-learn) (1.2.2)
Requirement already satisfied: scipy>=1.2.0 in /Users/divyabastola/opt/anaconda3/lib/python3.8/site-packages (from imbalanced-learn) (1.6.0)
Requirement already satisfied: numpy>=1.17.3 in /Users/divyabastola/opt/anaconda3/lib/python3.8/site-packages (from imbalanced-learn) (1.20.1)
Note: you may need to restart the kernel to use updated packages.
```

```
In [3]: conda install -c conda-forge imbalanced-learn

Collecting package metadata (current_repodata.json): done
Solving environment: done

==> WARNING: A newer version of conda exists. <==
  current version: 4.14.0
  latest version: 23.1.0

Please update conda by running

$ conda update -n base -c conda-forge conda

# All requested packages already installed.

Retrieving notices: ...working... done

Note: you may need to restart the kernel to use updated packages.
```

Import the CSV data set

```
In [4]: stroke_dataset = pd.read_csv('healthcare-dataset-stroke-data.csv')
stroke_dataset
stroke_dataset.describe()
```

	id	age	hypertension	heart_disease	avg_glucose_level	bmi	stroke
count	510.000000	510.000000	510.000000	510.000000	510.000000	4999.000000	510.000000
mean	3653.729254	43.226414	0.097458	0.084023	106.147677	28.969227	0.086728
std	2116.717252	12.416457	0.286671	0.236957	42.265500	7.954567	0.210230
min	0.00000000	0.000000	0.000000	0.000000	55.120000	10.300000	0.000000
25%	1774.125000	25.000000	0.000000	0.000000	77.245000	23.500000	0.000000
50%	3653.729254	43.000000	0.000000	0.000000	91.865000	28.100000	0.000000
75%	5456.200000	61.000000	0.000000	0.000000	114.000000	33.100000	0.000000
max	7294.000000	82.000000	1.000000	1.000000	271.740000	97.600000	1.000000

```
In [5]: stroke_dataset = pd.read_csv('healthcare-dataset-stroke-data.csv')
print(stroke_dataset.columns)

Index(['id', 'gender', 'age', 'hypertension', 'heart_disease', 'ever_married',
       'stroke_type', 'residence_type', 'avg_glucose_level', 'bmi',
       'smoking_status', 'stroke'],
      dtype='object')
```

Data Cleansing

The stroke_dataset is imported into the Jupyter along with all the libraries. For the first step of this project I have done data cleansing. Selevateral steps to data cleaning are:

1. First I have listed the data type of each column.
2. Secondly I checked if the data set had any null values.
3. In the third step, I handled missing values by deleting the records.
4. In the fourth step, I checked if there are any duplicate values especially in the id column.
5. I also deleted the column 'ever_married' due to multiple missing values.

```
In [6]: Step1 = stroke_dataset.drop('ever_married', axis=1)
print(Step1)

convert_dict = {'id': int,
                'gender': str,
                'age': int,
                'hypertension': int,
                'heart_disease': int,
                'residence_type': str,
                'avg_glucose_level': float,
                'bmi': float,
                'smoking_status': str,
                'stroke': int}

Step1 = Step1.astype(convert_dict)
print(Step1.dtypes)
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	0	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [7]: Step2 = smutils()
Step2.isnull().any()
Step2.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [8]: Stroke_df = Step2[Step2['bmi']!=0].dropna()
print(Stroke_df)
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [9]: Stroke_df.shape
(4999, 11)
```

```
In [10]: Step3 = smutils()
Step3.isnull().any()
Step3.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [11]: Step4 = smutils()
Step4.isnull().any()
Step4.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [12]: Step5 = smutils()
Step5.isnull().any()
Step5.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [13]: Step6 = smutils()
Step6.isnull().any()
Step6.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [14]: Step7 = smutils()
Step7.isnull().any()
Step7.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [15]: Step8 = smutils()
Step8.isnull().any()
Step8.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [16]: Step9 = smutils()
Step9.isnull().any()
Step9.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [17]: Step10 = smutils()
Step10.isnull().any()
Step10.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [18]: Step11 = smutils()
Step11.isnull().any()
Step11.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [19]: Step12 = smutils()
Step12.isnull().any()
Step12.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	
2	31112	Male	80	0	1	Private	
3	68182	Female	49	0	0	Private	
4	1665	Female	79	1	0	Self-employed	
5	56669	Male	81	0	0	Private	
...
5184	14180	Female	13	0	0	children	
5186	44873	Female	81	0	0	Self-employed	
5187	19723	Female	35	0	0	Self-employed	
5188	8269	Male	59	1	0	Private	
5189	44679	Female	44	0	0	Govt Job	

```
In [20]: Step13 = smutils()
Step13.isnull().any()
Step13.isnull().sum()
```

	id	gender	age	hypertension	heart_disease	work_type	
0	9046	Male	67	0	1	Private	
1	51076	Female	61	0	0	Self-employed	