

STROKE PREDICTION USING MACHINE LEARNING

In [1]:



```
1 import pandas as pd
2 import numpy as np
3 import warnings
4 warnings.filterwarnings('ignore')
5
6
7 # Data Visualization Libraries,
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import missingno as msv
11 from plotly.subplots import make_subplots
12 import plotly.graph_objects as go
13 from matplotlib.gridspec import GridSpec
14 from plotly.subplots import make_subplots
15 from plotly.offline import init_notebook_mode
16 from pywaffle import Waffle
17
18
19 # Machine Learning Libraries
20 from imblearn.over_sampling import SMOTE
21 from sklearn.metrics import precision_score, f1_score, recall_score, confusion_matrix
22 from sklearn.model_selection import train_test_split, RandomizedSearchCV, GridSearchCV
23 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
24 from sklearn.tree import DecisionTreeClassifier
25 from sklearn.linear_model import LogisticRegression
26 from sklearn.neighbors import KNeighborsClassifier
27
28 # Artificial Neural Network Libraries
29 import tensorflow as tf
30 from tensorflow import keras
31 from tensorflow.keras.utils import plot_model
32 from sklearn.preprocessing import StandardScaler
```

In [2]:



```
1 # Import your data
2 data = pd.read_csv('healthcare-dataset-stroke-data.csv')
```

In [3]:



```
1 #View the first 5 rows in the dataset
2 data.head()
```

Out[3]:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type
0	9046	Male	67.0	0	1	Yes	Private	Urban
1	51676	Female	61.0	0	0	Yes	Self-employed	Rural
2	31112	Male	80.0	0	1	Yes	Private	Rural
3	60182	Female	49.0	0	0	Yes	Private	Urban
4	1665	Female	79.0	1	0	Yes	Self-employed	Rural



In [4]:



```
1 # View Last 5 rows in the dataset
2 data.tail()
```

Out[4]:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_ty
5105	18234	Female	80.0	1	0	Yes	Private	Urb
5106	44873	Female	81.0	0	0	Yes	Self-employed	Urb
5107	19723	Female	35.0	0	0	Yes	Self-employed	Ru
5108	37544	Male	51.0	0	0	Yes	Private	Ru
5109	44679	Female	44.0	0	0	Yes	Govt_job	Urb



In [5]:



```
1 #get data description
2 data.describe()
```

Out[5]:

	id	age	hypertension	heart_disease	avg_glucose_level	
count	5110.000000	5110.000000	5110.000000	5110.000000	5110.000000	4909
mean	36517.829354	43.226614	0.097456	0.054012	106.147677	28
std	21161.721625	22.612647	0.296607	0.226063	45.283560	7
min	67.000000	0.080000	0.000000	0.000000	55.120000	10
25%	17741.250000	25.000000	0.000000	0.000000	77.245000	23
50%	36932.000000	45.000000	0.000000	0.000000	91.885000	28
75%	54682.000000	61.000000	0.000000	0.000000	114.090000	33
max	72940.000000	82.000000	1.000000	1.000000	271.740000	97



In [6]:



```
1 fig = make_subplots(rows=1, cols=2)
2
3 fig.add_trace(go.Indicator(
4     mode = "number",
5     value = data.shape[0],
6     number={'font':{'color': 'purple', 'size':75}},
7     delta = {"reference": 400},
8     title = {"text": "Rows <br><span style="
9         "'font-size:0.7em;color:gray'></span>"},
10    domain = {'y': [0, 1], 'x': [0, 0.2]}))
11 fig.add_trace(go.Indicator(
12     mode = "number",
13     value = data.shape[1],
14     number={'font':{'color': 'purple', 'size':75}},
15     delta = {"reference": 400},
16     title = {"text": "Columns <br><span style="
17         "'font-size:0.6em;color:gray'></span>"},
18    domain = {'y': [0, 0], 'x': [1, 1]}))
```

Rows

5110

Columns

12

In [6]:



```

1 # Created a function to get necessary information from the data set in a dataframe
2 def get_data_info():
3     dataset_info = pd.DataFrame(index=data.columns)
4     dataset_info['Data_type'] = data.dtypes
5     dataset_info['Total Value'] = data.count()
6     dataset_info['Null_count'] = data.isnull().sum()
7     dataset_info['Unique_count'] = data.nunique()
8     return dataset_info

```

In [7]:



```
1 get_data_info()
```

Out[7]:

	Data_type	Total Value	Null_count	Unique_count
id	int64	5110	0	5110
gender	object	5110	0	3
age	float64	5110	0	104
hypertension	int64	5110	0	2
heart_disease	int64	5110	0	2
ever_married	object	5110	0	2
work_type	object	5110	0	5
Residence_type	object	5110	0	2
avg_glucose_level	float64	5110	0	3979
bmi	float64	4909	201	418
smoking_status	object	5110	0	4
stroke	int64	5110	0	2

It is now easy to view missing data in all fields, their unique values and datatypes.

As seen above, the dataset only has missing values in the BMI column. Let's visualize that.

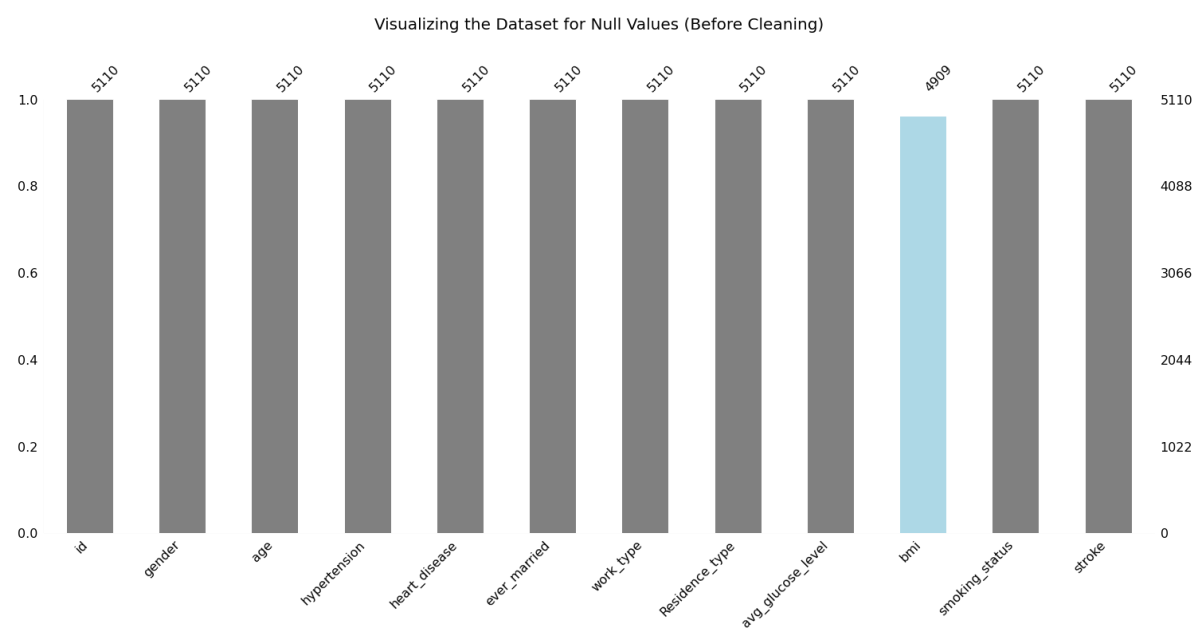
DATA CLEANING

In [8]:

```

1 MissingValuesColors = []
2
3 for i in data.columns:
4     if data[i].isna().sum() != 0:
5         MissingValuesColors.append('lightblue')
6     else:
7         MissingValuesColors.append('gray')
8
9 msv.bar(data, color=MissingValuesColors)
10 plt.title('Visualizing the Dataset for Null Values (Before Cleaning)', size=20, y=1.15)
11 plt.savefig("Visualizations\DataFrameDirty.png")
12 plt.show()

```



In [9]:

```
1 data.gender.value_counts()
```

Out[9]:

```

Female    2994
Male      2115
Other         1
Name: gender, dtype: int64

```

In [10]:



```
1 # Drop the others category because it is too small to consider.
```

In [11]:



```
1 #Drop row with others gender.  
2 data = data[data.gender != 'Other']
```

In [12]:



```
1 #Drop the missing values because the values are less than 5% of the total value  
2 data = data.dropna()
```

In [13]:



```
1 get_data_info()
```

Out[13]:

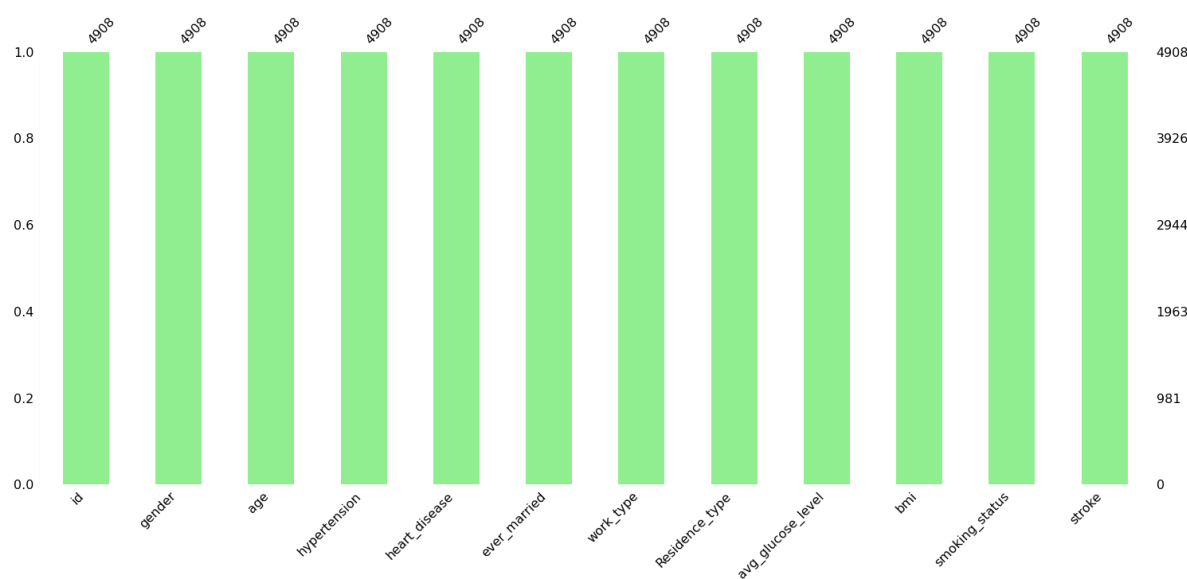
	Data_type	Total Value	Null_count	Unique_count
id	int64	4908	0	4908
gender	object	4908	0	2
age	float64	4908	0	104
hypertension	int64	4908	0	2
heart_disease	int64	4908	0	2
ever_married	object	4908	0	2
work_type	object	4908	0	5
Residence_type	object	4908	0	2
avg_glucose_level	float64	4908	0	3851
bmi	float64	4908	0	418
smoking_status	object	4908	0	4
stroke	int64	4908	0	2

In [14]:



```
1 MissingValuesColors = []
2
3 for i in data.columns:
4     if data[i].isna().sum() != 0:
5         MissingValuesColors.append('lightblue')
6     else:
7         MissingValuesColors.append('lightgreen')
8
9 msv.bar(data, color=MissingValuesColors)
10 plt.title('Visualizing the Dataset for Null Values (After Cleaning)', size=35, y=1.15)
11 plt.savefig("Visualizations\DataframeCleaned.png")
12 plt.show()
```

Visualizing the Dataset for Null Values (After Cleaning)



We have successfully removed all NA values.

In [15]:



```
1 # Age is seen as a float instead of Int, hence
2 # Convert Age to int
3 data['age'] = data['age'].astype(int)
4 get_data_info()
```

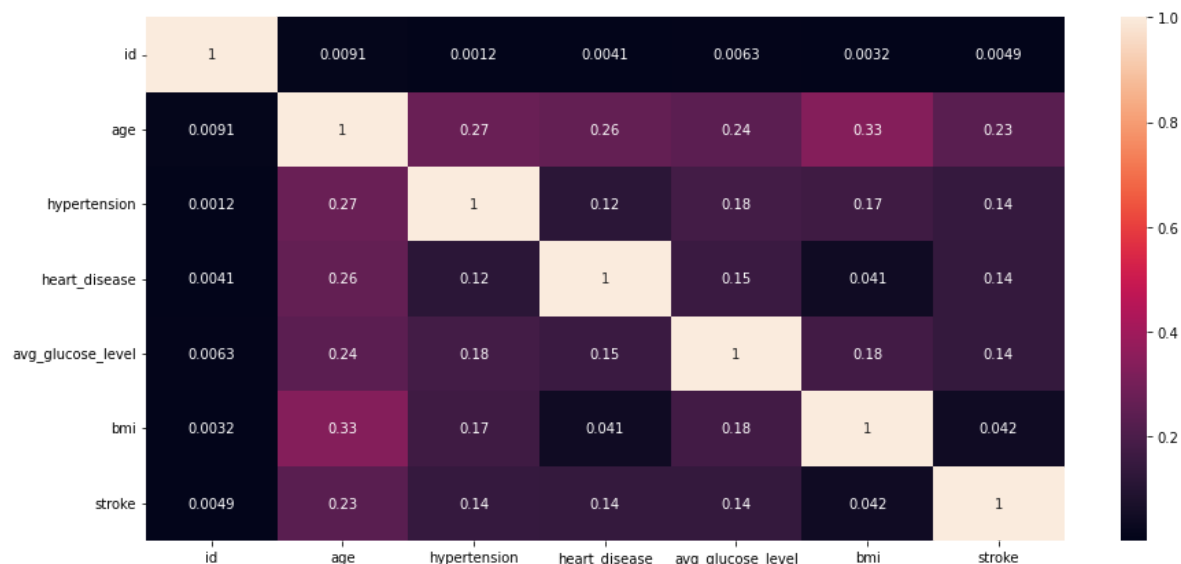
Out[15]:

	Data_type	Total Value	Null_count	Unique_count
id	int64	4908	0	4908
gender	object	4908	0	2
age	int32	4908	0	83
hypertension	int64	4908	0	2
heart_disease	int64	4908	0	2
ever_married	object	4908	0	2
work_type	object	4908	0	5
Residence_type	object	4908	0	2
avg_glucose_level	float64	4908	0	3851
bmi	float64	4908	0	418
smoking_status	object	4908	0	4
stroke	int64	4908	0	2

DATA VISUALIZATION

In [16]:

```
1 plt.figure(figsize=(15,7))
2 sns.heatmap(data.corr(),annot=True)
3 plt.savefig("Visualizations\heatmap.png")
4 plt.show()
```



In [17]:

```
1 #From the Dataset, Let's calculate the percentage of people that have stroke and do not
2 data.stroke.value_counts()
```

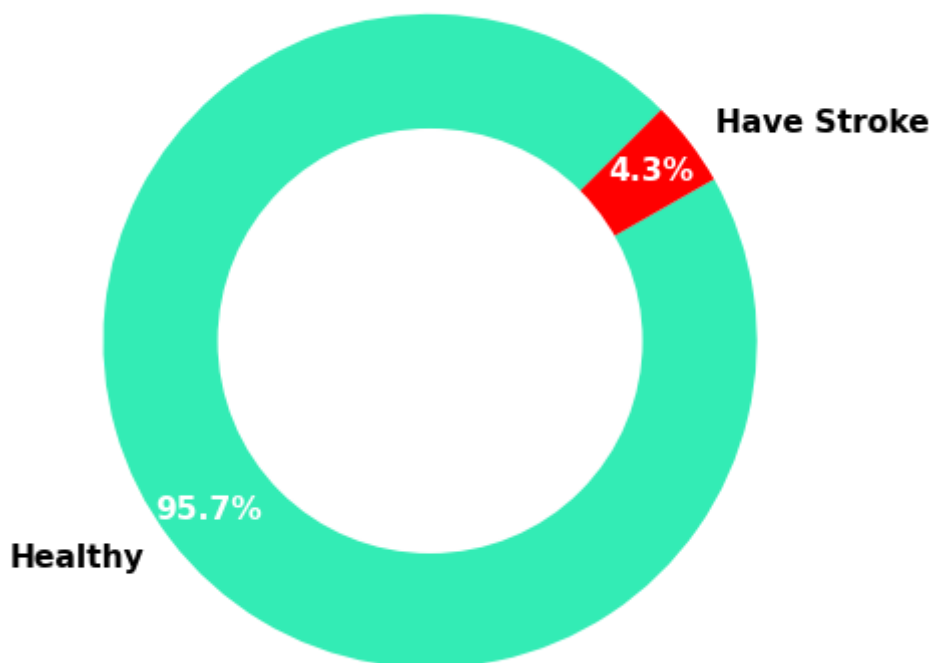
Out[17]:

```
0    4699
1     209
Name: stroke, dtype: int64
```

In [18]:

```
1 palette2 = ['#33ECB5', '#ff0000']
2
3 colors = ('#E2F11C', '#E3460A')
4 plt.figure(figsize=(10,7.5))
5 label = ['Healthy', 'Have Stroke']
6 patches, texts, pcts = plt.pie(data.stroke.value_counts(),
7                                labels=label,
8                                colors=[palette2[0], '#ff0000'],
9                                pctdistance=0.85,
10                               shadow=False,
11                               startangle=45,
12                               autopct='%1.1f%%',
13                               textprops={'fontsize': 15.5,
14                                         'weight': 'bold'})
15
16 plt.setp(pcts, color='white')
17
18 hfont = {'fontname': 'calibri', 'weight': 'bold'}
19 plt.title('Percentage of People living with Stroke', size=25, **hfont)
20
21 centre_circle = plt.Circle((0,0),0.65,fc='white')
22 fig = plt.gcf()
23 fig.gca().add_artist(centre_circle)
24 plt.savefig("Visualizations\PieChart.png")
25 plt.show()
```

Percentage of People living with Stroke



In [19]:



```
1 GenderGroupWithoutStroke = data.groupby(['gender', 'stroke']).count()['id'][[0,2]]
2 GenderGroupWithoutStroke
```

Out[19]:

```
gender  stroke
Female  0         2777
Male    0         1922
Name: id, dtype: int64
```

In [20]:



```
1 GenderGroupWithStroke = data.groupby(['gender', 'stroke']).count()['id'][[1,3]]
2 GenderGroupWithStroke
```

Out[20]:

```
gender  stroke
Female  1         120
Male    1          89
Name: id, dtype: int64
```

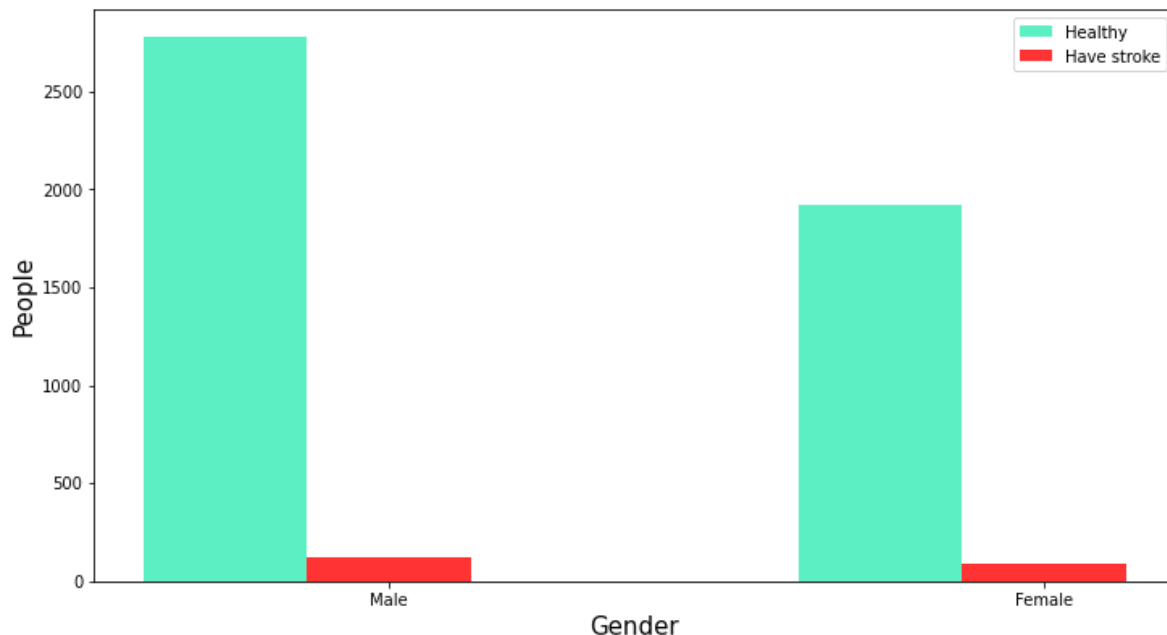
Assumption 1: MALE ARE MORE SUSCEPTIBLE TO STROKE DUE TO HIGH WORK RELATED STRESS

In [21]:



```
1 n_groups = 2
2
3 # create plot
4 plt.figure(figsize=(12,6.5))
5 indexx = np.arange(n_groups)
6 bar_width = 0.25
7 opacity = 0.8
8
9 rects1 = plt.bar(indexx, GenderGroupWithoutStroke, bar_width,
10 alpha=opacity,
11 color= palette2[0],
12 label='Healthy')
13
14 rects2 = plt.bar(indexx + bar_width, GenderGroupWithStroke, bar_width,
15 alpha=opacity,
16 color='#ff0000',
17 label='Have stroke')
18
19 plt.xlabel('Gender', size=15)
20 plt.ylabel('People',size=15)
21 plt.title('Gender influence on stroke\n',size=25, **hfont)
22 plt.xticks(indexx + bar_width, ('Male', 'Female'))
23 plt.legend()
24 plt.savefig("Visualizations\GenderInfluence.png")
25 plt.show()
```

Gender influence on stroke



In [22]:

```

1  #Let's visualize the data by finding the percentage of people living with stroke and w
2  stroke_gen = data[data['stroke'] == 1]['gender'].value_counts()
3  healthy_gen = data[data['stroke'] == 0]['gender'].value_counts()
4
5  female = data['gender'].value_counts().values[0]
6  male = data['gender'].value_counts().values[1]
7
8  stroke_female = int(round(stroke_gen.values[0] / female * 100, 0))
9  stroke_male = int(round(stroke_gen.values[1] / male * 100, 0))
10 healthy_female = int(round(healthy_gen.values[0] / female * 100, 0))
11 healthy_male = int(round(healthy_gen.values[1] / male * 100, 0))
12
13 female_per = int(round(female/(female+male) * 100, 0))
14 male_per = int(round(male/(female+male) * 100, 0))
15
16 fig = plt.figure(FigureClass = Waffle,
17                  constrained_layout = True,
18                  figsize = (7,7),
19                  facecolor = '#fff',dpi = 100,
20
21                  plots = {'121':
22                          {
23                              'rows':6,
24                              'columns': 6,
25                              'values' : [healthy_male,stroke_male],
26                              'colors' : [palette2[0], '#ff0000'],
27                              'vertical' : True,
28                              'interval_ratio_y': 0.1,
29                              'interval_ratio_x': 0.1,
30                              'icons' : 'male',
31                              'icon_legend': False,
32                              'icon_size':20,
33                              'plot_anchor':'C',
34                              'alpha':0.1
35                          },
36
37                          '122' :
38                          {
39                              'rows': 6,
40                              'columns':6,
41                              'values':[healthy_female,stroke_female],
42                              'colors' : [palette2[0], '#ff0000'],
43                              'vertical': True,
44                              'interval_ratio_y': 0.1,
45                              'interval_ratio_x': 0.1,
46                              'icons' : 'female',
47                              'icon_legend' :False,
48                              'icon_size':20,
49                              'plot_anchor':'C',
50                              'alpha':0.1
51                          }
52                      },
53
54
55 )
56 fig.text(0., 0.8, "Gender that's more susceptible to stroke", {'font':'Serif', 'size':1
57 fig.text(0.23, 0.28, '{}%'.format(healthy_male), {'font':'Serif', 'size':20,'weight':'
58 fig.text(0.65, 0.28, '{}%'.format(healthy_female), {'font':'Serif', 'size':20,'weight'
59 fig.text(0.6,0.73, 'Stroke ', {'font': 'Serif','weight':'bold','Size': '16','weight':'

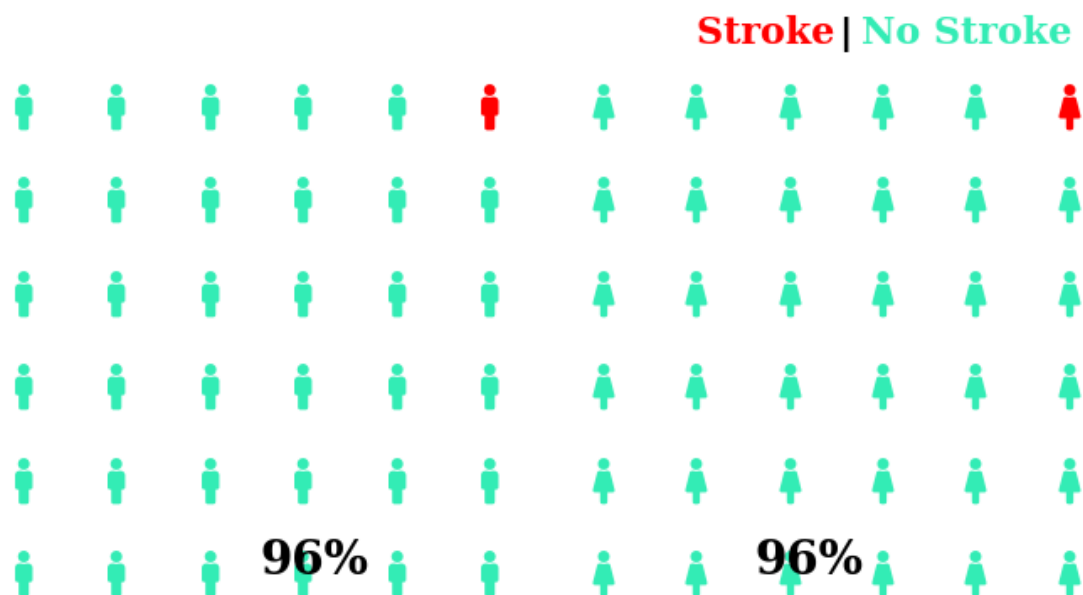
```

```

60 fig.text(0.72,0.73, '|', {'color':'black' , 'size':'16', 'weight': 'bold'})
61 fig.text(0.74,0.73, 'No Stroke', {'font': 'Serif','weight':'bold', 'Size': '16','style
62 plt.savefig("Visualizations\GenderProportion.png")
63 fig.show()

```

Gender that's more susceptible to stroke



i From the above diagram, it shows that the assumption is wrong. Being a male doesn't mean you're more susceptible to stroke.

? Question 1) DOES AGE HAVE AN IMPACT ON STROKE?

In [23]:

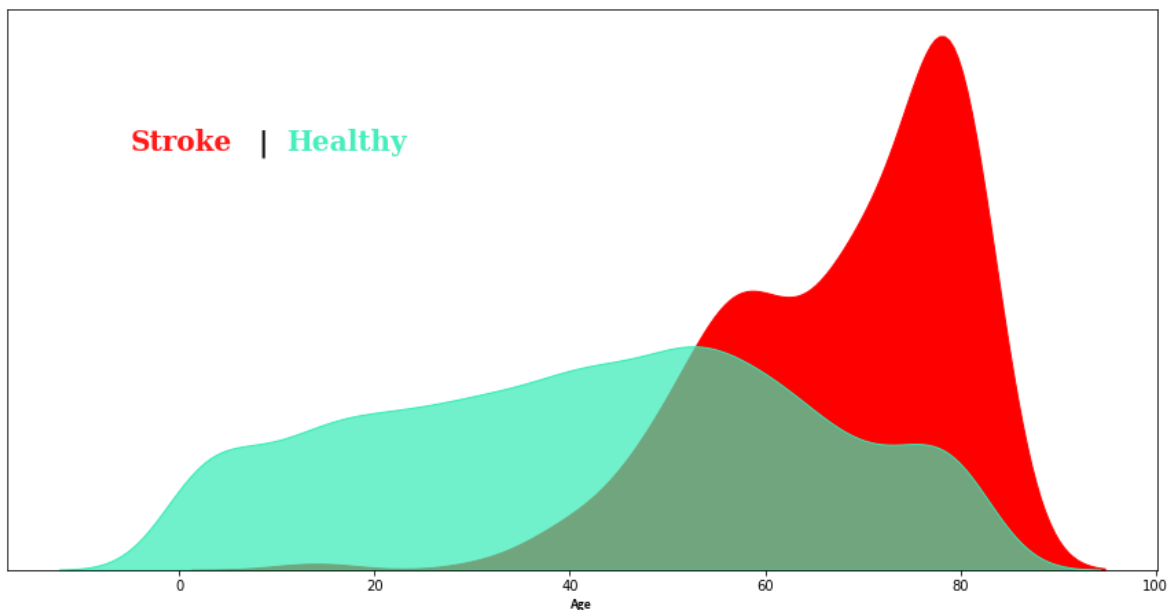


```

1 fig = plt.figure(figsize=(15, 7.5))
2 ax = fig.add_subplot(111)
3 plt.title('Distribution of People by Age\n', size=28, **hfont)
4 ax.grid(False)
5 ax.axes.get_yaxis().set_visible(False)
6 ax.text(-5, 0.03, 'Stroke', {'font': 'Serif',
7                               'size': '20',
8                               'weight': 'bold',
9                               'color': '#ff0000'}, alpha=0.9)
10
11 ax.text(8, 0.03, '|', {'font': 'Serif',
12                        'size': '20',
13                        'weight': 'bold',
14                        'color': 'black'}, alpha=0.9)
15
16 ax.text(11, 0.03, 'Healthy', {'font': 'Serif',
17                               'size': '20',
18                               'weight': 'bold',
19                               'color': palette2[0]}, alpha=0.9)
20
21 sns.kdeplot(data=data[data.stroke == 1],
22             x='age', shade=True, ax=ax, color='#ff0000', alpha=1)
23 sns.kdeplot(data=data[data.stroke == 0],
24             x='age', shade=True, ax=ax, color=palette2[0], alpha=0.7)
25 plt.xlabel('Age', **hfont)
26 plt.savefig("Visualizations\DistributionOfPeopleByAge.png")
27 plt.show()

```

Distribution of People by Age



In [24]:



```
1 #Categorizing my BMI, Age and Glucose Columns into categorical values using range.
2 #Categories used
3 #Age: 0-12 : Children,      BMI: 0-18 : Underweight      GlucoseLevel: 0-89
4 #   13-17 : Teens,          19-24 : Ideal          90-159
5 #   18-44 : Adults          25-30 : Overweight      160-229
6 #   45-59 : Mid Adults      30-50 : Obesity          230-500
7 #   60-150 : Elderly
8
9 data['bmi_cat'] = pd.cut(data['bmi'], bins = [0, 19, 25,30,50], labels = ['Underweight
10 data['age_cat'] = pd.cut(data['age'], bins = [0,13,18, 45,60,150], labels = ['Children
11 data['glucose_cat'] = pd.cut(data['avg_glucose_level'], bins = [0,90,160,230,500], label
```



In [25]:



```

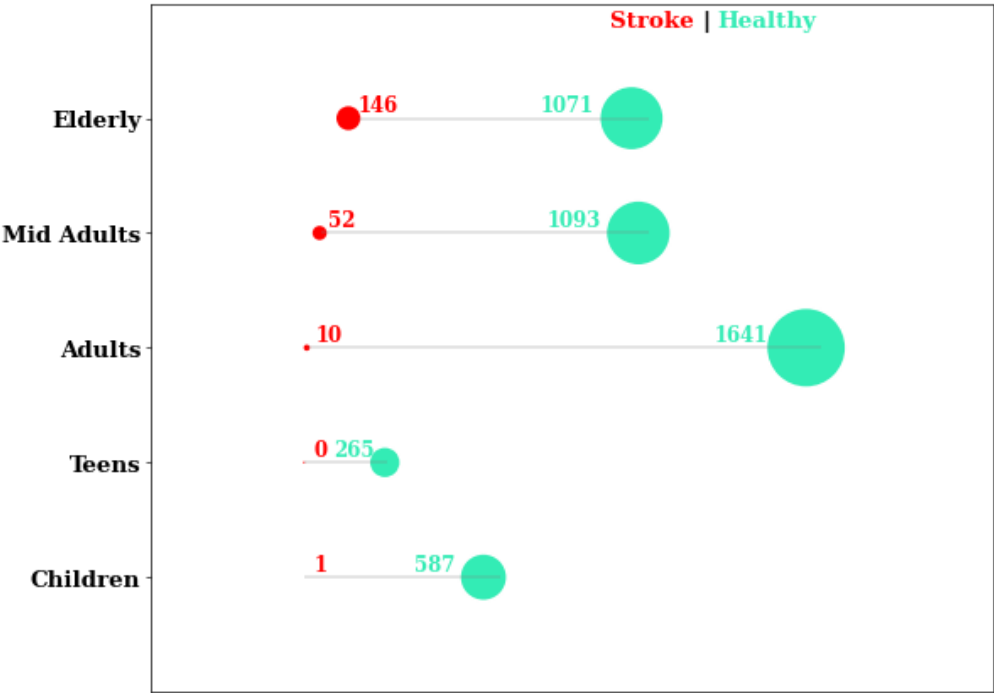
1  # Using the Categories created above, Let's Visualize those that have Stroke with the A
2
3  fig = plt.figure(figsize = (24,10), dpi = 60)
4
5  gs = fig.add_gridspec(10,24)
6  gs.update(wspace = 1, hspace = 0.05)
7
8  ax1 = fig.add_subplot(gs[1:10,13:]) #dumbbell plot
9  stroke_age = data[data['stroke'] == 1].age_cat.value_counts()
10 healthy_age = data[data['stroke'] == 0].age_cat.value_counts()
11
12 ax1.hlines(y = ['Children', 'Teens', 'Adults', 'Mid Adults', 'Elderly'], xmin = [644,2
13           xmax = [1,1,11,59,177], color = 'grey',**{'linewidth':0.5})
14
15
16 sns.scatterplot(y = stroke_age.index, x = stroke_age.values, s = stroke_age.values*2, c
17 sns.scatterplot(y = healthy_age.index, x = healthy_age.values, s = healthy_age.values*2
18
19 ax1.axes.get_xaxis().set_visible(False)
20 ax1.set_xlim(xmin = -500, xmax = 2250)
21 ax1.set_ylim(ymin = -1,ymax = 5)
22
23 ax1.set_yticklabels( labels = ['Children', 'Teens', 'Adults', 'Mid Adults', 'Elderly']).
24
25 ax1.text(0,5.8, 'Impact of age on stroke? \n' ,{'font': 'Serif', 'Size': '20','weight'
26 ax1.text(1000,4.8, '\nStroke ', {'font': 'Serif','weight':'bold','Size': '16','weight'
27 ax1.text(1300,4.8, '\n|', {'color':'black' , 'size':'16', 'weight': 'bold'})
28 ax1.text(1350,4.8, '\nHealthy', {'font': 'Serif','weight':'bold', 'Size': '16','style'
29 ax1.text(-550,5., 'Age has a significant impact on stroke, and clearly seen that stroke
30     {'font':'Serif', 'size':'16','color': 'black'})
31
32 ax1.text(stroke_age.values[0] + 30,4.05, stroke_age.values[0], {'font':'Serif', 'Size'
33 ax1.text(healthy_age.values[2] - 300,4.05, healthy_age.values[2], {'font':'Serif', 'Si
34
35 ax1.text(stroke_age.values[1] + 30,3.05, stroke_age.values[1], {'font':'Serif', 'Size'
36 ax1.text(healthy_age.values[1] - 300,3.05, healthy_age.values[1], {'font':'Serif', 'Si
37
38 ax1.text(stroke_age.values[2] + 30,2.05, stroke_age.values[2], {'font':'Serif', 'Size'
39 ax1.text(healthy_age.values[0] - 300,2.05, healthy_age.values[0], {'font':'Serif', 'Si
40
41 ax1.text(stroke_age.values[4] + 30,1.05, stroke_age.values[4], {'font':'Serif', 'Size'
42 ax1.text(healthy_age.values[4] - 170,1.05, healthy_age.values[4], {'font':'Serif', 'Si
43
44 ax1.text(stroke_age.values[3] + 30,0.05, stroke_age.values[3], {'font':'Serif', 'Size'
45 ax1.text(healthy_age.values[3] - 225,0.05, healthy_age.values[3], {'font':'Serif', 'Si
46 plt.savefig("Visualizations\ImpactofAgeonStroke.png")
47 plt.show()

```



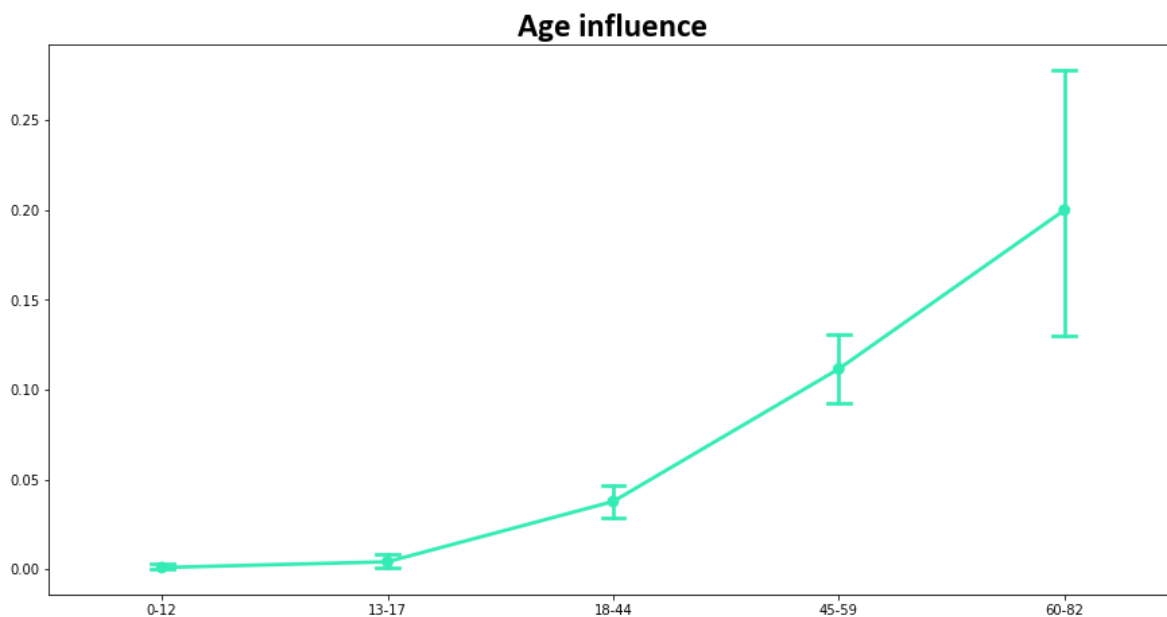
Impact of age on stroke?

Age has a significant impact on stroke, and clearly seen that stroke numbers are highest for elderly people and mid age adults, where as negligible for younger people.



In [26]:

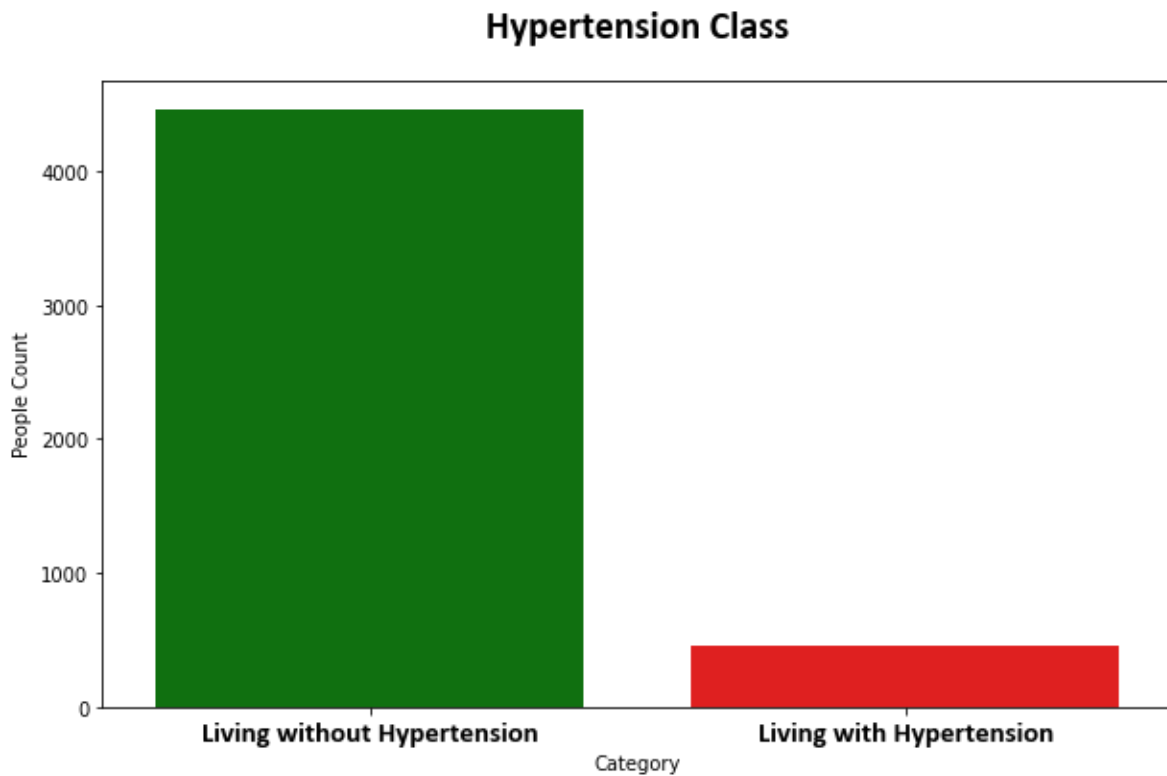
```
1 label = ['0-12', '13-17', '18-44', '45-59', f'60-{round(data.age.max())}']
2
3 def plot_age(data, Column_Name):
4
5     AgeData = data[[Column_Name, 'stroke']]
6     AgeData[Column_Name] = pd.cut(AgeData[Column_Name],
7                                   bins=[0, 20, 40, 60, 80, 100],
8                                   labels=label)
9
10    color = np.random.choice(palette2, 1)[0]
11    plt.figure(figsize=(15, 7.5))
12    plot = sns.pointplot(x=Column_Name, y='stroke',
13                        dodge=0.1, capsize=.1, data=AgeData, color=color)
14    plot.set_title(f'Age influence', fontsize=25, **hfont)
15    plot.set(xlabel=None, ylabel=None)
16    plt.savefig("Visualizations\AgeInfluence.png")
17    plt.show()
18 plot_age(data, 'age')
```



In [27]:



```
1 plt.figure(figsize=(10, 6))
2 HypertensiveDistribution = sns.countplot(x=data.hypertension,palette=(['green','red'])).
3 HypertensiveDistribution.set_title('Hypertension Class',fontsize=20, y=1.05, **hfont)
4 HypertensiveDistribution.set(xlabel='Category', ylabel='People Count')
5 plt.grid(False)
6 HypertensiveDistribution.set_xticklabels(['Living without Hypertension', 'Living with H
7 plt.savefig("Visualizations\HypertensionClass.png")
8 plt.show()
```



In [28]:



```
1 data.smoking_status.unique()
```

Out[28]:

```
array(['formerly smoked', 'never smoked', 'smokes', 'Unknown'],  
      dtype=object)
```

In [29]:



```
1 data.columns
```

Out[29]:

```
Index(['id', 'gender', 'age', 'hypertension', 'heart_disease', 'ever_married',  
      'work_type', 'Residence_type', 'avg_glucose_level', 'bmi',  
      'smoking_status', 'stroke', 'bmi_cat', 'age_cat', 'glucose_cat'],  
      dtype='object')
```

In [30]:



```
1 HypertensionGroupWithoutStroke = data.groupby(['hypertension', 'stroke']).count()['id']  
2 HypertensionGroupWithoutStroke
```

Out[30]:

```
hypertension  stroke  
0            0      4308  
            1      149  
Name: id, dtype: int64
```

In [31]:



```
1 HypertensionGroupWithStroke = data.groupby(['hypertension', 'stroke']).count()['id']  
2 HypertensionGroupWithStroke
```

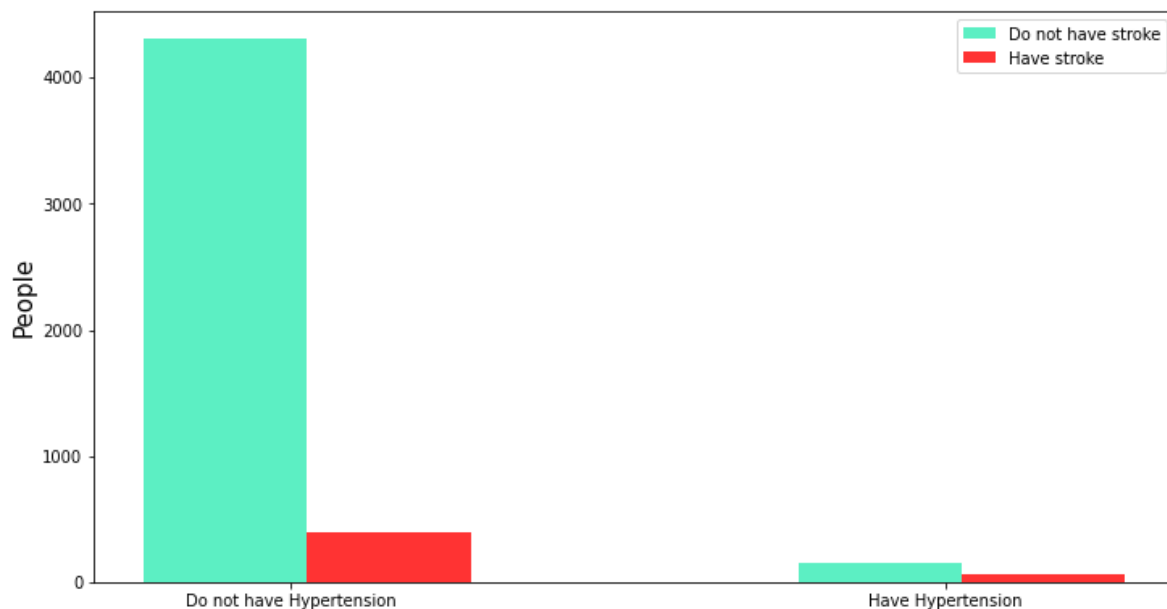
Out[31]:

```
hypertension  stroke  
1            0      391  
            1       60  
Name: id, dtype: int64
```

In [32]:

```
1 n_groups = 2
2
3 # create plot
4 plt.figure(figsize=(12,6.5))
5 indexx = np.arange(n_groups)
6 bar_width = 0.25
7 opacity = 0.8
8
9 rects1 = plt.bar(indexx, HypertensionGroupWithoutStroke, bar_width,
10 alpha=opacity,
11 color= palette2[0],
12 label='Do not have stroke')
13
14 rects2 = plt.bar(indexx + bar_width, HypertensionGroupWithStroke, bar_width,
15 alpha=opacity,
16 color='#ff0000',
17 label='Have stroke')
18
19 plt.xlabel('', size=15)
20 plt.ylabel('People',size=15)
21 plt.title('Hypertension influence on stroke\n',size=25, **hfont)
22 plt.xticks(indexx + 0.1,('Do not have Hypertension', 'Have Hypertension'))
23 plt.legend()
24 plt.savefig("Visualizations\HypertensionInfluenceOnStroke.png")
25 plt.show()
```

Hypertension influence on stroke



From the above diagram, it shows that there is an higher risk for people that have Hypertension to have stroke

In [33]:

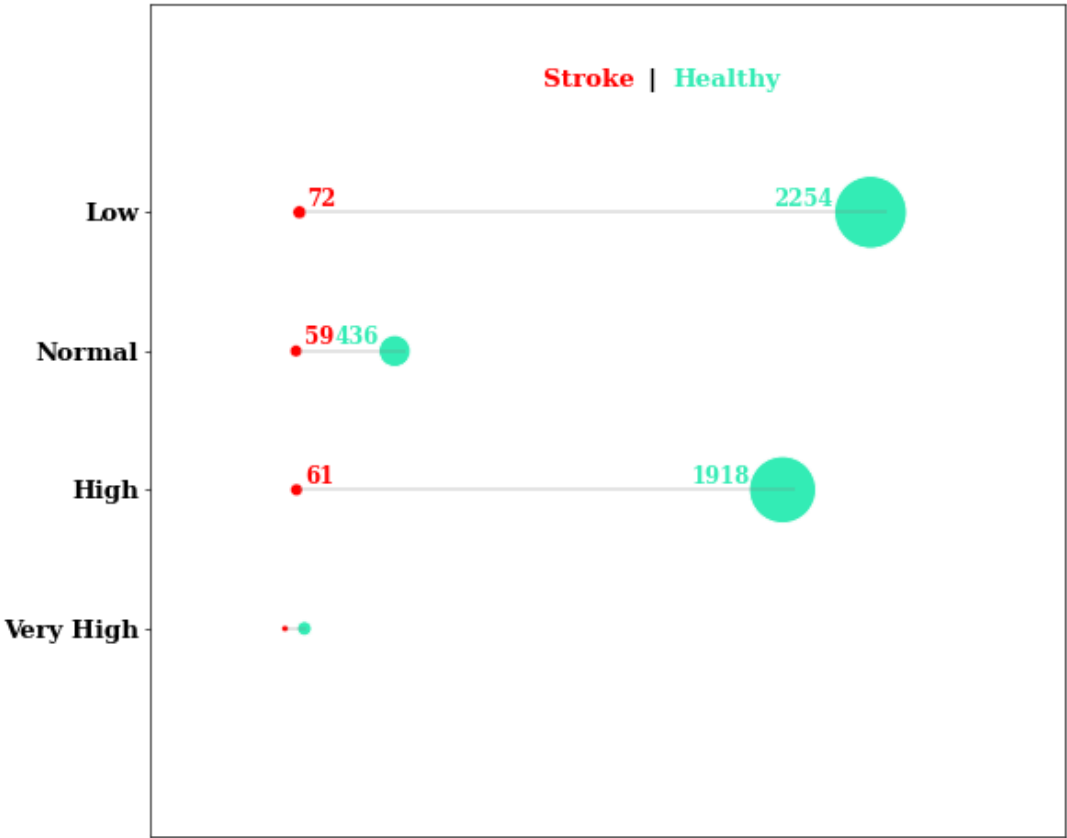
```

1  # Let's Visualize those that have Stroke using the Glucose Level Category using a dumbbell
2
3  fig = plt.figure(figsize = (24,10), dpi = 60)
4
5  gs = fig.add_gridspec(10,24)
6  gs.update(wspace = 1, hspace = 0.05)
7
8  ax1 = fig.add_subplot(gs[0:,13:]) #dumbbell plot
9  stroke_glucose = data[data['stroke'] == 1].glucose_cat.value_counts()
10 healthy_glucose = data[data['stroke'] == 0].glucose_cat.value_counts()
11
12 ax1.hlines(y = ['Very High', 'High', 'Normal', 'Low'], xmin = [18,71,71,89],
13           xmax = [101,1966,478,2316], color = 'grey',**{'linewidth':0.5})
14
15
16 sns.scatterplot(y = stroke_glucose.index, x = stroke_glucose.values[[0,2,1,3]], s = stroke_glucose.values[[0,2,1,3]])
17 sns.scatterplot(y = healthy_glucose.index, x = healthy_glucose.values[[0,2,1,3]], s = healthy_glucose.values[[0,2,1,3]])
18
19 ax1.axes.get_xaxis().set_visible(False)
20 ax1.set_xlim(xmin = -500, xmax = 3000)
21 ax1.set_ylim(ymin = -1.5,ymax = 4.5)
22
23 ax1.set_yticklabels( labels = ['Very High', 'High', 'Normal', 'Low'],fontdict = {'font'
24
25 ax1.text(0,5.8, 'How Glucose level impacts Having Stroke' ,{'font': 'Serif', 'Size': '16'})
26 ax1.text(1000,3.9, '\nStroke ', {'font': 'Serif','weight':'bold','Size': '16','weight'
27 ax1.text(1400,3.9, '\n|', {'color':'black' , 'size':'16', 'weight': 'bold'})
28 ax1.text(1500,3.9, '\nHealthy', {'font': 'Serif','weight':'bold', 'Size': '16','style'
29 ax1.text(120,5., 'Glucose does not have a significant impact on strokes\n',
30         {'font':'Serif', 'size':'16','color': 'black'})
31
32 ax1.text(stroke_glucose.values[0] + 30,3.05, stroke_glucose.values[0], {'font':'Serif',
33 ax1.text(healthy_glucose.values[0] - 370,3.05, healthy_glucose.values[0], {'font':'Serif',
34
35 ax1.text(stroke_glucose.values[2] + 30,2.05, stroke_glucose.values[2], {'font':'Serif',
36 ax1.text(healthy_glucose.values[2] - 230,2.05, healthy_glucose.values[2], {'font':'Serif',
37
38 ax1.text(stroke_glucose.values[1] + 30,1.05, stroke_glucose.values[1], {'font':'Serif',
39 ax1.text(healthy_glucose.values[1] - 350,1.05, healthy_glucose.values[1], {'font':'Serif',
40 plt.savefig("Visualizations\GlucoseLevelImpactonStroke.png")
41 plt.show()

```

How Glucose level impacts Having Stroke

Glucose does not have a significant impact on strokes



In [34]:

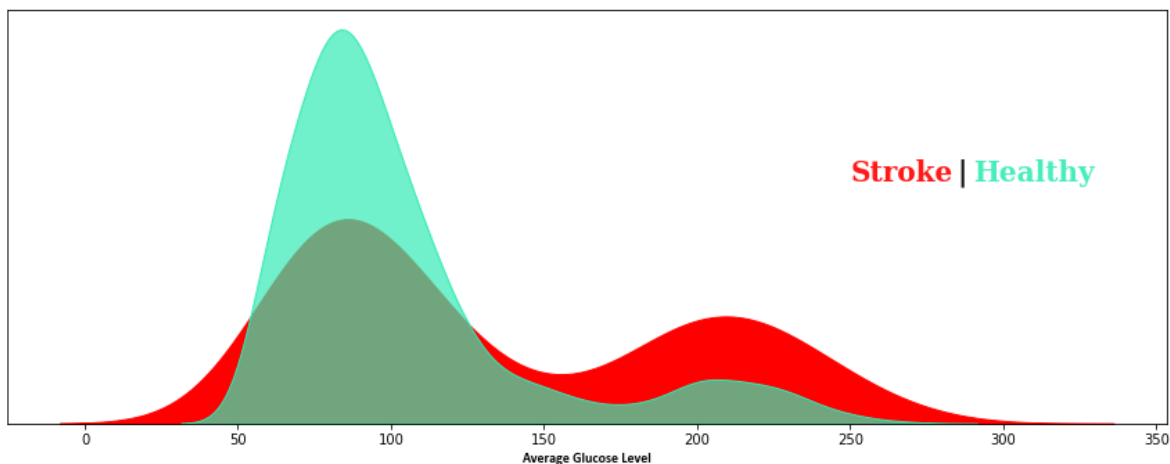


```

1 fig = plt.figure(figsize=(15, 5.5))
2 ax = fig.add_subplot(111)
3 plt.title('Glucose-Stroke Distribution of People\n', size=28, **hfont)
4 ax.grid(False)
5 ax.axes.get_yaxis().set_visible(False)
6 ax.text(250, 0.01, 'Stroke', {'font': 'Serif',
7                               'size': '20',
8                               'weight': 'bold',
9                               'color': '#ff0000'}, alpha=0.9)
10
11 ax.text(285, 0.01, '|', {'font': 'Serif',
12                           'size': '20',
13                           'weight': 'bold',
14                           'color': 'black'}, alpha=0.9)
15
16 ax.text(290, 0.01, 'Healthy', {'font': 'Serif',
17                                 'size': '20',
18                                 'weight': 'bold',
19                                 'color': palette2[0]}, alpha=0.9)
20
21 sns.kdeplot(data=data[data.stroke == 1],
22             x='avg_glucose_level', shade=True, ax=ax, color='#ff0000', alpha=1)
23 sns.kdeplot(data=data[data.stroke == 0],
24             x='avg_glucose_level', shade=True, ax=ax, color=palette2[0], alpha=0.7)
25 plt.xlabel('Average Glucose Level', **hfont)
26 plt.savefig("Visualizations\GlucoseLevelonStroke.png")
27 plt.show()

```

Glucose-Stroke Distribution of People



In [35]:

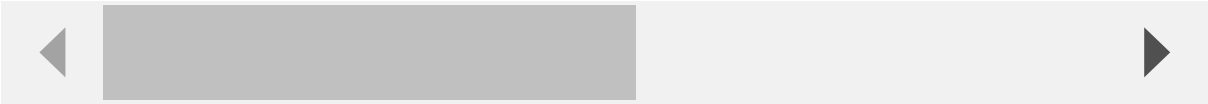


```
1 data[data.stroke == 1]
```

Out[35]:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type
0	9046	Male	67	0	1	Yes	Private	Urban
2	31112	Male	80	0	1	Yes	Private	Rural
3	60182	Female	49	0	0	Yes	Private	Urban
4	1665	Female	79	1	0	Yes	Self-employed	Rural
5	56669	Male	81	0	0	Yes	Private	Urban
...
243	40460	Female	68	1	1	Yes	Private	Urban
244	17739	Male	57	0	0	Yes	Private	Rural
245	49669	Female	14	0	0	No	children	Rural
246	27153	Female	75	0	0	Yes	Self-employed	Rural
248	43424	Female	78	0	0	Yes	Private	Rural

209 rows × 15 columns



In [36]:



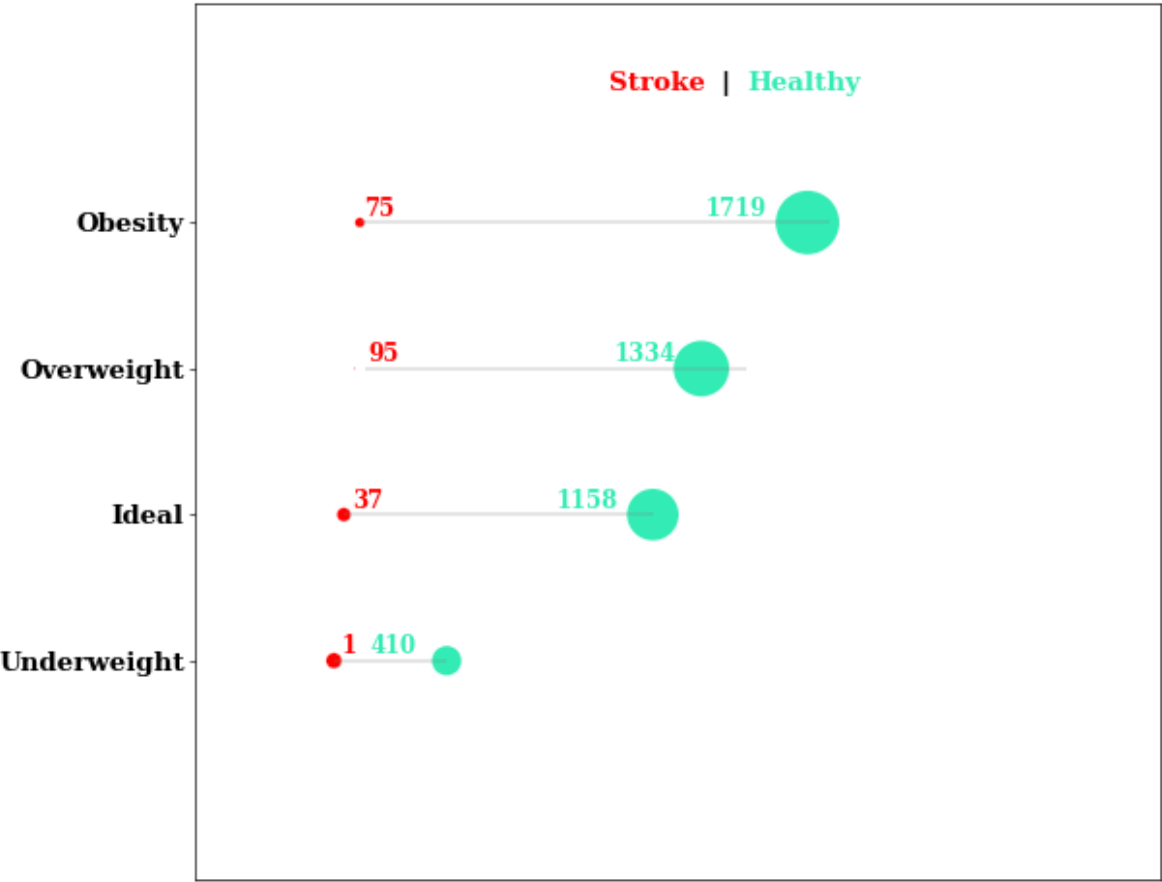
```

1  # Let's Visualize those that have Stroke using the BMI Levels Category using a dumbbell
2
3  fig = plt.figure(figsize = (24,10), dpi = 60)
4
5  gs = fig.add_gridspec(10,24)
6  gs.update(wspace = 1, hspace = 0.05)
7
8  ax1 = fig.add_subplot(gs[0:,13:]) #dumbbell plot
9  stroke_bmi = data[data['stroke'] == 1].bmi_cat.value_counts()
10 healthy_bmi = data[data['stroke'] == 0].bmi_cat.value_counts()
11
12 ax1.hlines(y = ['Underweight','Ideal','Overweight','Obesity'], xmin = [1,37,115,96],
13           xmax = [410,1158,1495,1797], color = 'grey',**{'linewidth':0.5})
14
15
16 sns.scatterplot(y = stroke_bmi.index[[3,2,0,1]], x = stroke_bmi.values[[3,2,0,1]], s =
17 sns.scatterplot(y = healthy_bmi.index[[0,1,2,3]], x = healthy_bmi.values[[0,1,2,3]], s =
18
19 ax1.axes.get_xaxis().set_visible(False)
20 ax1.set_xlim(xmin = -500, xmax = 3000)
21 ax1.set_ylim(ymin = -1.5,ymax = 4.5)
22
23 ax1.set_yticklabels( labels = ['Underweight','Ideal','Overweight','Obesity'],fontdict :
24
25 ax1.text(0,5.8, 'How BMI impacts on Having Stroke ',{'font': 'Serif', 'Size': '20','we:
26 ax1.text(1000,3.9, '\nStroke ', {'font': 'Serif','weight':'bold','Size': '16','weight'
27 ax1.text(1400,3.9, '\n|', {'color':'black' , 'size':'16', 'weight': 'bold'})
28 ax1.text(1500,3.9, '\nHealthy', {'font': 'Serif','weight':'bold', 'Size': '16','style'
29 ax1.text(120,5., 'BMI has a significant impact on stroke\n',
30         {'font':'Serif', 'size':'16','color': 'black'})
31
32 ax1.text(stroke_bmi.values[1] + 30,3.05, stroke_bmi.values[1], {'font':'Serif', 'Size'
33 ax1.text(healthy_bmi.values[0] - 370,3.05, healthy_bmi.values[0], {'font':'Serif', 'Si:
34
35 ax1.text(stroke_bmi.values[0] + 30,2.05, stroke_bmi.values[0], {'font':'Serif', 'Size'
36 ax1.text(healthy_bmi.values[1] - 320,2.05, healthy_bmi.values[1], {'font':'Serif', 'Si:
37
38 ax1.text(stroke_bmi.values[2] + 30,1.05, stroke_bmi.values[2], {'font':'Serif', 'Size'
39 ax1.text(healthy_bmi.values[2] - 350,1.05, healthy_bmi.values[2], {'font':'Serif', 'Si:
40
41 ax1.text(stroke_bmi.values[3] + 30,0.05, stroke_bmi.values[3], {'font':'Serif', 'Size'
42 ax1.text(healthy_bmi.values[3] - 280,0.05, healthy_bmi.values[3], {'font':'Serif', 'Si:
43 plt.savefig("Visualizations\BMILevelOnStroke.png")
44 plt.show()
45

```

How BMI impacts on Having Stroke

BMI has a significant impact on stroke



In [37]:



```

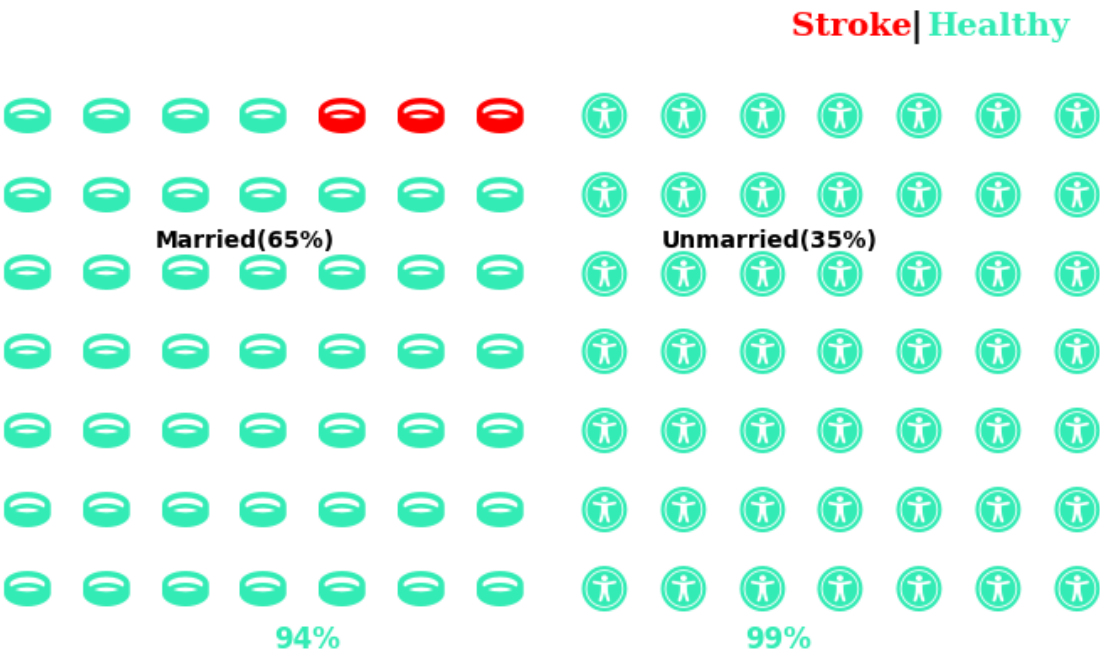
1 stroke_married = data[data['stroke'] == 1]['ever_married'].value_counts()
2 healthy_married = data[data['stroke'] == 0]['ever_married'].value_counts()
3
4 yes = data['ever_married'].value_counts().values[0]
5 no = data['ever_married'].value_counts().values[1]
6
7 stroke_no = int(round(stroke_married.values[1] / no * 100, 0))
8 stroke_yes = int(round(stroke_married.values[0] / yes * 100, 0))
9 healthy_no = int(round(healthy_married.values[1] / no * 100, 0))
10 healthy_yes = int(round(healthy_married.values[0] / yes * 100, 0))
11
12 no_per = int(round(no/(no+yes) * 100, 0))
13 yes_per = int(round(yes/(no+yes) * 100, 0))
14
15
16 fig = plt.figure(FigureClass = Waffle,
17                 constrained_layout = True,
18                 figsize = (7,7),
19                 facecolor = '#fff',dpi = 100,
20
21                 plots = {'121':
22                     {
23                         'rows':7,
24                         'columns': 7,
25                         'values' : [stroke_yes,healthy_yes],
26                         'colors' : ['#ff0000',palette2[0]],
27                         'vertical' : True,
28                         'interval_ratio_x': 0.005,
29                         'interval_ratio_y': 0.005,
30                         'icons' : 'ring',
31                         'icon_legend': False,
32                         'icon_size':20,
33                         'plot_anchor':'C',
34                         'alpha':0.8,
35                         'starting_location': 'NE'
36                     },
37
38                     '122' :
39                     {
40                         'rows': 7,
41                         'columns':7,
42                         'values':[stroke_no,healthy_no],
43                         'colors' : ['#ff0000',palette2[0]],
44                         'vertical': True,
45                         'interval_ratio_x': 0.005,
46                         'interval_ratio_y':0.005,
47                         'icons' : 'universal-access',
48                         'icon_legend' :False,
49                         'icon_size':20,
50                         'plot_anchor':'C',
51                         'alpha':0.8,
52                         'starting_location': 'NE'
53                     }
54                 },
55
56 )
57
58 fig.text(0., 0.8, 'Does being married have a deciding effect on people living with stroke?')
59 fig.text(0.25, 0.23, '{}%'.format(healthy_yes), {'size':12,'weight':'bold' , 'color':palette2[0]})

```


5/19/22, 4:54 PMB1232321_Oladeinbo_Olutayo_Notebook - Jupyter Notebook

```
60 fig.text(0.65, 0.23, '{}%'.format(healthy_no), {'size':12,'weight':'bold', 'color':'palegreen'})
61 fig.text(0.15, 0.57, 'Married({}% )'.format(yes_per), {'size':10,'weight':'bold', 'color':'black'})
62 fig.text(0.58, 0.57, "Unmarried({}%)".format(no_per), {'size':10,'weight':'bold', 'color':'black'})
63 fig.text(0.69,0.75, 'Stroke ', {'font': 'Serif','weight':'bold','Size': '14','weight':'bold'})
64 fig.text(0.79,0.75, '|', {'color':'black' , 'size':'14', 'weight': 'bold'})
65 fig.text(0.805,0.75, 'Healthy', {'font': 'Serif','weight':'bold', 'Size': '14','style':'italic'})
66 plt.savefig("Visualizations\MarriageDistribution.png")
67 fig.show()
```

Does being married have a deciding effect on people living with stroke?



Risk of stroke in married people is relatively high

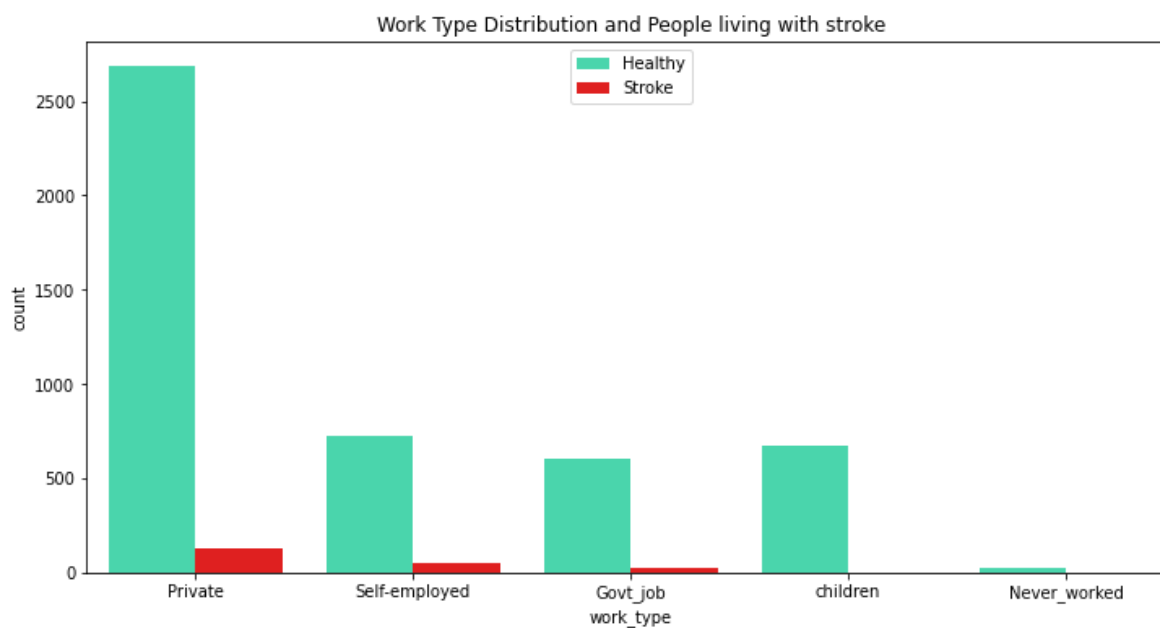
In [38]:

```
1 # Work Type
```

In [39]:



```
1 # using countplot to plot for Work Type
2 plt.figure(figsize = (12,6))
3 labels=['Healthy','Stroke']
4 s = sns.countplot(x="work_type", hue="stroke", data=data,palette=palette2)
5 h,l = s.get_legend_handles_labels()
6 s.legend(h,labels, title="")
7 plt.title("Work Type Distribution and People living with stroke",size=12)
8 plt.savefig("Visualizations\WorkTypeDistribution.png")
9 plt.show()
```



In [40]:

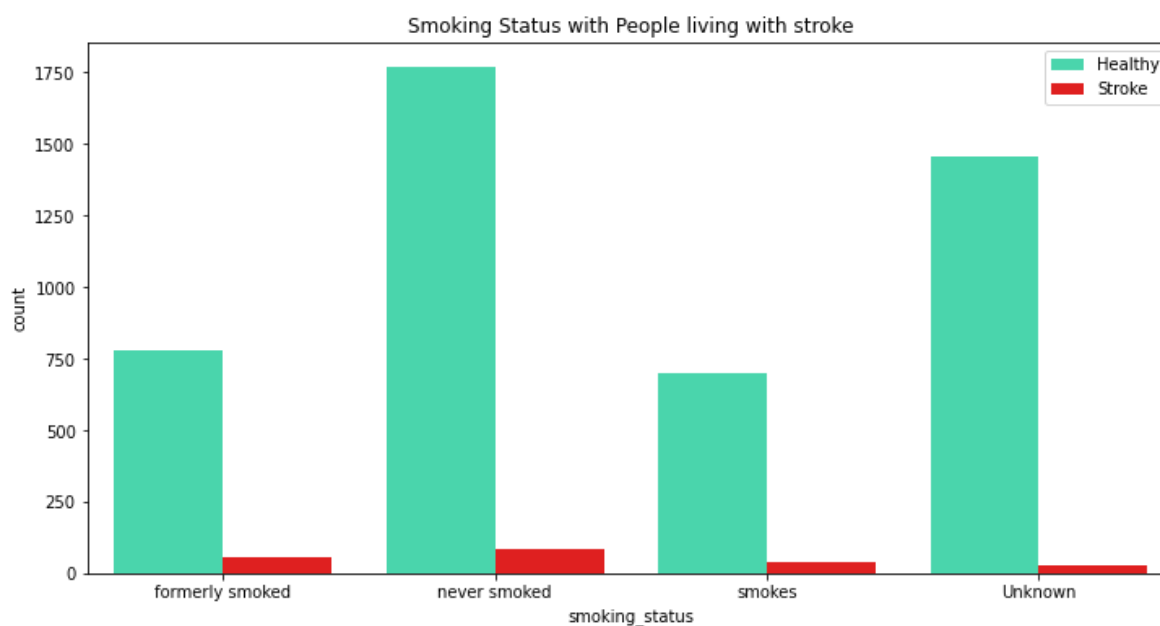


```
1 # Smoking
```

In [41]:



```
1 # using countplot to plot for Work Type
2 plt.figure(figsize = (12,6))
3 labels=['Healthy','Stroke']
4 s = sns.countplot(x="smoking_status", hue="stroke", data=data,palette=palette2)
5 h,l = s.get_legend_handles_labels()
6 s.legend(h,labels, title="")
7 plt.title("Smoking Status with People living with stroke",size=12)
8 plt.savefig("Visualizations\SmokingStatus.png")
9 plt.show()
```



Modelling and Results

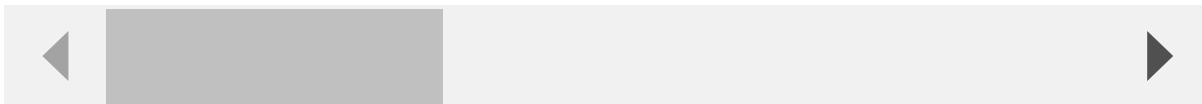
In [44]:



```
1 data = data.drop(['bmi_cat', 'age_cat', 'glucose_cat'],axis='columns')
2 data = pd.get_dummies(data,drop_first=True)
3 data.head()
```

Out[44]:

	id	age	hypertension	heart_disease	avg_glucose_level	bmi	stroke	gender_Male	ever
0	9046	67	0	1	228.69	36.6	1	1	
2	31112	80	0	1	105.92	32.5	1	1	
3	60182	49	0	0	171.23	34.4	1	0	
4	1665	79	1	0	174.12	24.0	1	0	
5	56669	81	0	0	186.21	29.0	1	1	



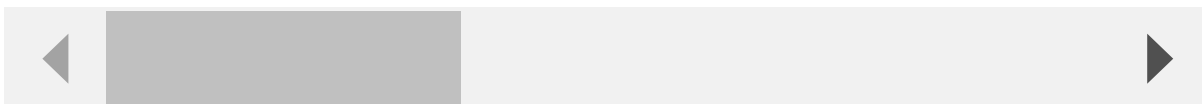
In [45]:



```
1 X = data.drop(['stroke', 'id'],axis='columns')
2 X.head()
```

Out[45]:

	age	hypertension	heart_disease	avg_glucose_level	bmi	gender_Male	ever_married_Yes
0	67	0	1	228.69	36.6	1	1
2	80	0	1	105.92	32.5	1	1
3	49	0	0	171.23	34.4	0	1
4	79	1	0	174.12	24.0	0	1
5	81	0	0	186.21	29.0	1	1



In [46]:



```
1 y = data.stroke
```

In [47]:



```
1 X.shape,y.shape
```

Out[47]:

((4908, 15), (4908,))

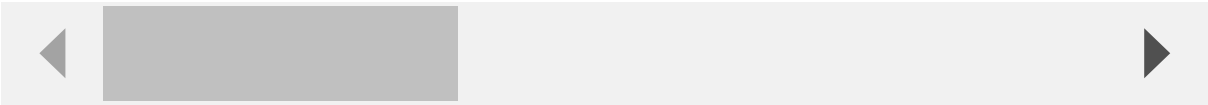
In [48]:



```
1 X.head()
```

Out[48]:

	age	hypertension	heart_disease	avg_glucose_level	bmi	gender_Male	ever_married_Yes
0	67	0	1	228.69	36.6	1	1
2	80	0	1	105.92	32.5	1	1
3	49	0	0	171.23	34.4	0	1
4	79	1	0	174.12	24.0	0	1
5	81	0	0	186.21	29.0	1	1



In [49]:



```
1 y[:5]
```

Out[49]:

0 1
2 1
3 1
4 1
5 1
Name: stroke, dtype: int64

HYPERPARAMETER TUNING

In [46]:



```
1 #Python Dictionary with 5 supervised models and parameters to choose the best Models ar
2
3 model_params = {
4     'Decision Tree': {
5         'model' : DecisionTreeClassifier(),
6         'params' : {
7             'criterion':['gini','entropy'],
8             'splitter': ['best','random'],
9             'max_depth': [10,20,30,100],
10            'random_state': [1,2,10]
11        }
12    },
13    'Random_forest':{
14        'model' : RandomForestClassifier(),
15        'params' : {
16            'n_estimators': [1,5,100],
17            'n_jobs': [1,10,20],
18            'random_state': [1,2,10]
19        }
20    },
21    'Logistic_regression' :{
22        'model' : LogisticRegression(),
23        'params': {
24            'C': [1,5,10],
25            'solver':['liblinear','saga'],
26            'multi_class':['auto'],
27            'random_state': [1,2,10],
28            'penalty': ['l1','l2','elasticnet','none']
29        }
30    },
31    'K_Nearest_Neighbour' :{
32        'model' : KNeighborsClassifier(),
33        'params' :{
34            'n_neighbors': [1,5,10],
35            'algorithm': ["auto", "brute", "kd_tree", "ball_tree"],
36            'weights': ['uniform','distance'],
37            'n_jobs' : [1,10,20]
38        }
39    },
40    'Gradient_Boost': {
41        'model': GradientBoostingClassifier(),
42        'params' :{
43            'learning_rate': [0.01],
44            'loss': ['exponential'],
45            'max_depth': [50,70],
46            'max_features': [1,2],
47            'n_estimators': [200,300]
48        }
49    }
50 }
```

In [47]:

```

1 scores = [] #check list comprehension
2
3 for model_name,mp in model_params.items():
4     X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=1,
5     sm = SMOTE(random_state=0)
6     X_train,y_train = sm.fit_resample(X_train,y_train)
7     scaler = StandardScaler()
8     X_train = scaler.fit_transform(X_train)
9     rs = GridSearchCV(mp['model'],mp['params'],cv=5,return_train_score=False)
10    rs.fit(X_train,y_train)
11    scores.append({
12        'Model': model_name,
13        'Best_Score': rs.best_score_,
14        'Best_Parameters':rs.best_params_
15    })

```

In [48]:

```

1 pd.options.display.max_colwidth = 200
2 scoresdf = pd.DataFrame(scores,columns=['Model','Best_Score','Best_Parameters'])
3 scoresdf.sort_values(by='Best_Score',ascending=False, inplace=True)
4 scoresdf

```

Out[48]:

	Model	Best_Score	Best_Parameters
4	Gradient_Boost	0.956374	{'learning_rate': 0.01, 'loss': 'exponential', 'max_depth': 70, 'max_features': 2, 'n_estimators': 300}
1	Random_forest	0.952251	{'n_estimators': 100, 'n_jobs': 1, 'random_state': 10}
3	K_Nearest_Neighbour	0.929505	{'algorithm': 'auto', 'n_jobs': 1, 'n_neighbors': 1, 'weights': 'uniform'}
0	Decision Tree	0.927110	{'criterion': 'gini', 'max_depth': 30, 'random_state': 2, 'splitter': 'random'}
2	Logistic_regression	0.858744	{'C': 1, 'multi_class': 'auto', 'penalty': 'l2', 'random_state': 1, 'solver': 'liblinear'}

In [49]:

```

1 # The above table shows the Model and Parameters with the best score in a descending order
2 # So we use the Models and the Parameters displayed in our table above for all our algorithms

```

In [50]:

```

1 # Function that accepts (X,y,model_name,model,random_state,Datasplit and parameters per model)
2 # the model score and stores it in model_results.

```

In [50]:



```

1 models_results = {}
2
3 def show_model_results(X,y,model_name,model,rand_state,Datasplit=0.2,**kwargs):
4     print(f'The model {model_name} with parameters : {kwargs}')
5     # Create an object m of the model with parameters entered into the function
6     m = model(**kwargs)
7     # Split data into training and testing set
8     X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=Datasplit,random_state=
9     # Create an object of a SMOTE (Oversampling Library)
10    sm = SMOTE(random_state=0)
11    # Performing oversampling on our train set
12    X_train,y_train = sm.fit_resample(X_train,y_train)
13    # Create an object of our scaling class
14    scaler = StandardScaler()
15    # Scale our X set
16    X_train = scaler.fit_transform(X_train)
17    X_test = scaler.transform(X_test)
18    # Model training
19    m.fit(X_train,y_train)
20    # Get model score
21    score = m.score(X_test,y_test)
22    # Get model predictions
23    prediction = m.predict(X_test)
24    # Get model precision score
25    model_precision = precision_score(y_test,prediction)
26    # Get model recall score
27    recall = recall_score(y_test,prediction)
28    # Get model F1 score
29    F1 = f1_score(y_test,prediction)
30    print('')
31    print('*****')
32    print(f'Model name:         \t {model_name}')
33    print(f'Model Parameters:  \t {kwargs}')
34    print(f'Model Score:         \t {score*100}')
35    print(f'Model Precision Score: {model_precision*100}')
36    print(f'Model Recall Score:  \t {recall*100}')
37    print(f'Model F1 Score:   \t {F1*100}')
38    print('*****')
39    # Call function plot_confusion matrix
40    plot_confusion_matrices(m,X_test,y_test,model_name)
41    return score,model_name,F1,model_precision,recall
42 # Function to plot our models confusion matrix
43 def plot_confusion_matrices(model, X_test, y_test,model_name):
44     # Get model prediction
45     y_pred = model.predict(X_test)
46     # confusion matrix
47     matrix = confusion_matrix(y_test, y_pred)
48     # Dataframe to store values
49     df_cm = pd.DataFrame(matrix, index = ['Stroke', 'Healthy'],
50                             columns = ['Stroke', 'Healthy'])
51     plt.figure(figsize = (12,8))
52     #plot confusion matrix
53     sns.heatmap(df_cm,
54                 annot=True,
55                 cmap='Greens',
56                 fmt='.5g',
57                 annot_kws={"size": 20}).set_title('Confusion matrix', fontsize = 35, y:
58     plt.xlabel('Predicted values', fontsize = 20, **hfont)

```



```
59 plt.ylabel('True values', fontsize = 20, **hfont)
60 plt.savefig(f"Visualizations\{model_name}.png")
61 plt.show()
```

1. RANDOM FOREST CLASSIFIER

In [51]:



```

1 rnd_state = 2
2 # Call show model result and pass parameters from Hyperparameter tuning
3 output = show_model_results(X,y,'Random Forest',RandomForestClassifier,n_estimators=100)
4 # Save outputs in a dictionary
5 models_resultsRC = ({
6     'Model Name': output[1],
7     'Model Score': output[0],
8     'F1 Score': output[2],
9     'Precision Score': output[3],
10    'Recall Score': output[4]
11 })

```

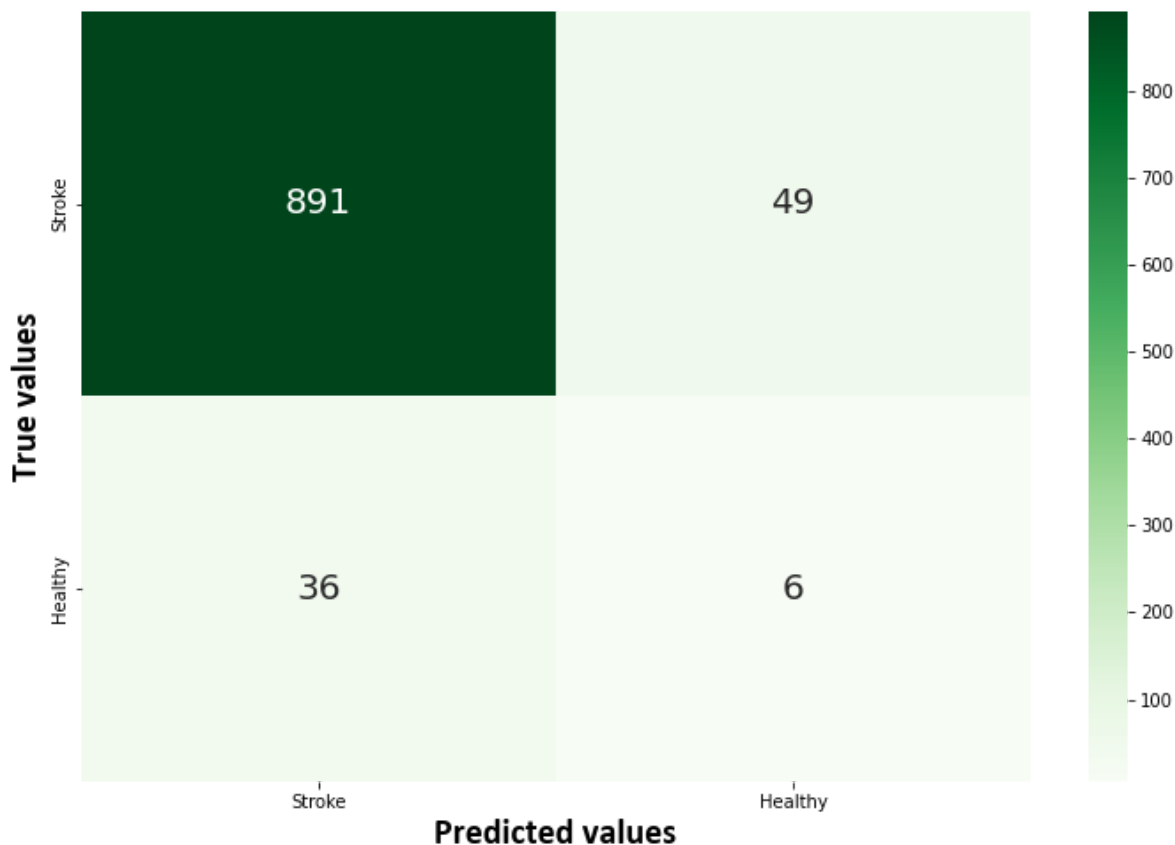
The model Random Forest with parameters : {'n_estimators': 100, 'n_jobs': 1}

```

*****
Model name:          Random Forest
Model Parameters:    {'n_estimators': 100, 'n_jobs': 1}
Model Score:         91.34419551934828
Model Precision Score: 10.909090909090908
Model Recall Score:  14.285714285714285
Model F1 Score:      12.371134020618557
*****

```

Confusion matrix



2. K NEAREST NEIGHBOURS

In [52]:



```

1 rnd_state = 1
2 # Call show model result and pass parameters from Hyperparameter tuning
3 output = show_model_results(X,y,'KNN',KNeighborsClassifier,algorithm='auto',n_jobs=1,n
4 # Save outputs in a dictionary
5 models_resultsKN = ({
6     'Model Name': output[1],
7     'Model Score': output[0],
8     'F1 Score': output[2],
9     'Precision Score': output[3],
10    'Recall Score': output[4]
11 })

```

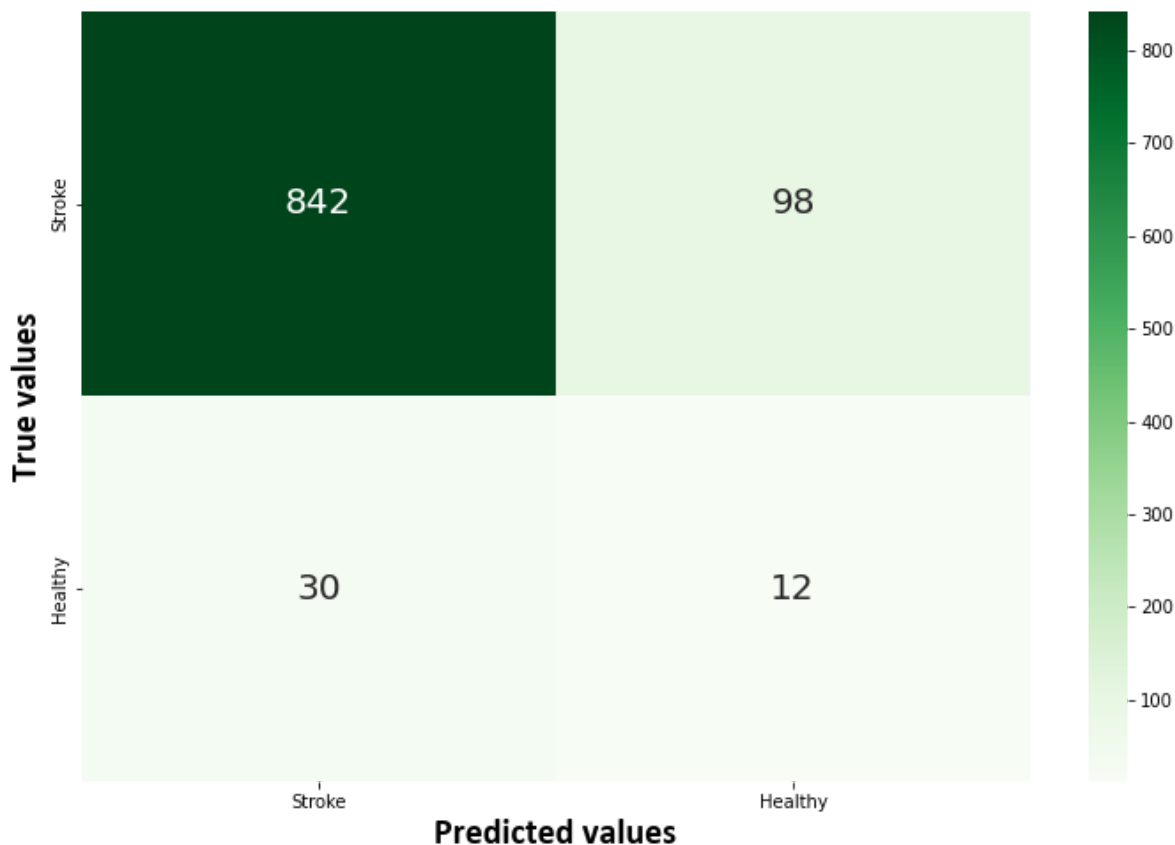
The model KNN with parameters : {'algorithm': 'auto', 'n_jobs': 1, 'n_neighbors': 1, 'weights': 'uniform'}

```

*****
Model name:                KNN
Model Parameters:          {'algorithm': 'auto', 'n_jobs': 1, 'n_neighbors':
1, 'weights': 'uniform'}
Model Score:                86.9653767820774
Model Precision Score:      10.909090909090908
Model Recall Score:         28.57142857142857
Model F1 Score:             15.789473684210526
*****

```

Confusion matrix



3. Gradient Boost

In [53]:



```

1 rand_state = 0
2 # Call show model result and pass parameters from Hyperparameter tuning
3 output = show_model_results(X, y, 'Gradient Boost', GradientBoostingClassifier, learning_
4 # Save outputs in a dictionary
5 models_resultsGB = ({
6     'Model Name': output[1],
7     'Model Score': output[0],
8     'F1 Score': output[2],
9     'Precision Score': output[3],
10    'Recall Score': output[4]
11 })

```

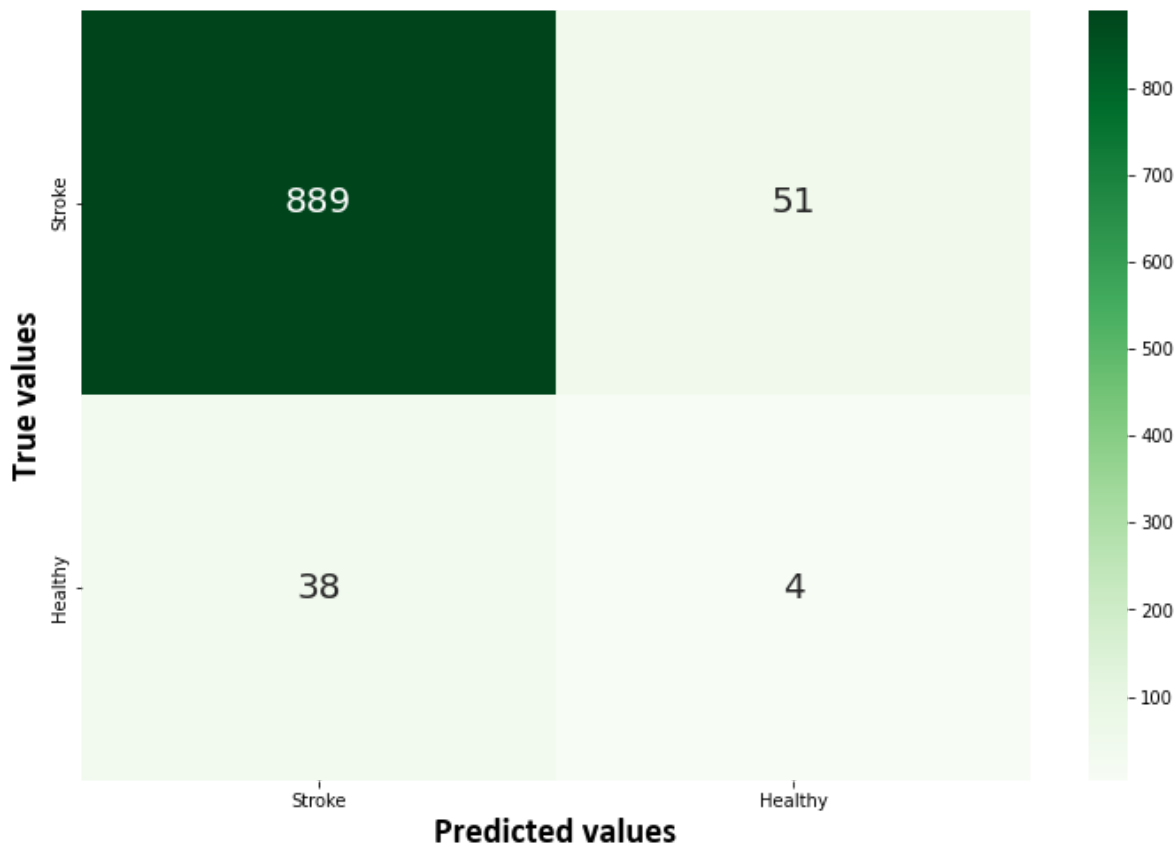
The model Gradient Boost with parameters : {'learning_rate': 0.01, 'loss': 'exponential', 'max_depth': 70, 'max_features': 1, 'n_estimators': 200}

```

Model name:          Gradient Boost
Model Parameters:    {'learning_rate': 0.01, 'loss': 'exponential', 'max_
_depth': 70, 'max_features': 1, 'n_estimators': 200}
Model Score:         90.93686354378818
Model Precision Score: 7.27272727272725
Model Recall Score:  9.523809523809524
Model F1 Score:      8.24742268041237

```

Confusion matrix



4. Decision Tree

In [54]:



```

1 # Call show model result and pass parameters from Hyperparameter tuning
2 output = show_model_results(X,y,'Decision Tree',DecisionTreeClassifier,10,criterion='entropy')
3 # Save outputs in a dictionary
4 models_resultsDT = ({
5     'Model Name': output[1],
6     'Model Score': output[0],
7     'F1 Score': output[2],
8     'Precision Score': output[3],
9     'Recall Score': output[4]
10 })

```



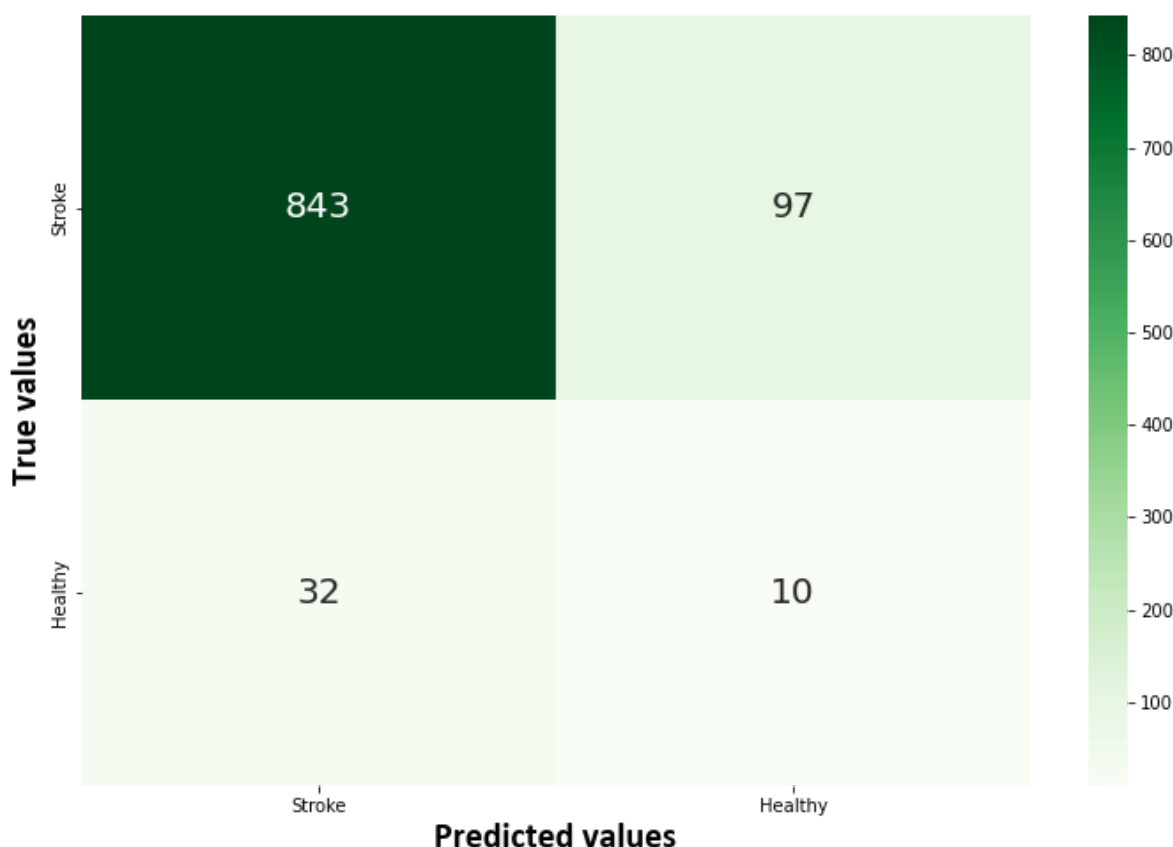
The model Decision Tree with parameters : {'criterion': 'entropy', 'max_depth': 20, 'random_state': 2, 'splitter': 'random'}

```

Model name:          Decision Tree
Model Parameters:    {'criterion': 'entropy', 'max_depth': 20, 'random_state': 2, 'splitter': 'random'}
Model Score:         86.86354378818737
Model Precision Score: 9.345794392523365
Model Recall Score:  23.809523809523807
Model F1 Score:      13.422818791946309

```

Confusion matrix



5. Logistic Regression

In [55]:



```

1 # Call show model result and pass parameters from Hyperparameter tuning
2 output = show_model_results(X,y,'Logistic Regression',LogisticRegression,2,C=1,multi_c
3 # Save outputs in a dictionary
4 models_resultsLR = ({
5     'Model Name': output[1],
6     'Model Score': output[0],
7     'F1 Score': output[2],
8     'Precision Score': output[3],
9     'Recall Score': output[4]
10 })

```

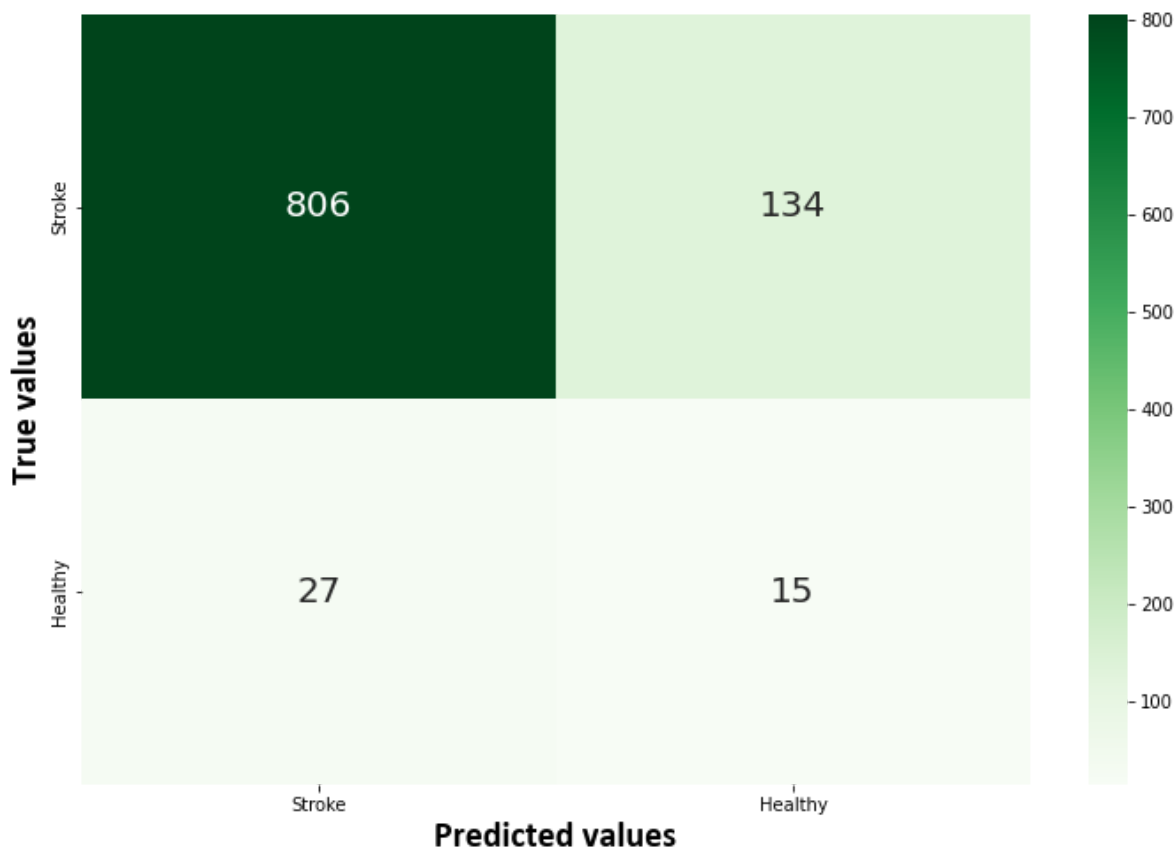
The model Logistic Regression with parameters : {'C': 1, 'multi_class': 'auto', 'penalty': 'none', 'solver': 'saga'}

```

*****
Model name:          Logistic Regression
Model Parameters:    {'C': 1, 'multi_class': 'auto', 'penalty': 'none',
'solver': 'saga'}
Model Score:         83.60488798370672
Model Precision Score: 10.06711409395973
Model Recall Score:   35.714285714285715
Model F1 Score:       15.706806282722512
*****

```

Confusion matrix



6. Artificial Neural Network

In [168]:

```
1 sm = SMOTE(random_state=0)
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2, shuffle=True)
4 X_train, y_train = sm.fit_resample(X_train, y_train)
5 y_train = keras.utils.to_categorical(y_train, 2)
6 y_test = keras.utils.to_categorical(y_test, 2)
7 scaler = StandardScaler()
8 X_train = scaler.fit_transform(X_train)
9 X_test = scaler.transform(X_test)
10 # Creating a Neural Network with 1 input Layer and 3 hidden layers with activation function relu
11 # Then one output layer with ***sigmoid*** function
12 model = keras.Sequential([
13     keras.layers.Flatten(input_dim=X_train.shape[1]),
14     keras.layers.Dense(500, activation='relu'),
15     keras.layers.Dense(250, activation='relu'),
16     keras.layers.Dense(125, activation='relu'),
17     keras.layers.Dense(2, activation='sigmoid')
18 ])
19 #Compile our model using Optimizer adam and Loss function ,
20 model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
21 # Model trains for 150 epochs and validates our model on X_test and y_test.
22 history = model.fit(X_train, y_train, epochs=150, validation_data=(X_test, y_test), verbose=1)
```

In [169]:

```
1 y_predicted = model.predict(X_test)
2 y_predicted
```

Out[169]:

```
array([[9.99879479e-01, 5.42283058e-04],
       [1.00000000e+00, 2.25540669e-10],
       [5.42799234e-02, 4.18436527e-01],
       ...,
       [9.96947706e-01, 5.72471619e-02],
       [1.00000000e+00, 1.26689295e-17],
       [1.00000000e+00, 3.40698151e-29]], dtype=float32)
```

In [170]:

```
1 y_predicted_labels = [np.argmax(i) for i in y_predicted]
2 y_predicted_labels[:5]
```

Out[170]:

```
[0, 0, 1, 0, 0]
```

In [171]:



```
1 y_test_labels = [np.argmax(i) for i in y_test]
2 y_test_labels[:5]
```

Out[171]:

```
[0, 0, 0, 0, 0]
```

In [173]:



```
1 from sklearn.metrics import classification_report
2 # Print-Out our classification Report
3 print(classification_report(y_test_labels,y_predicted_labels))
```

	precision	recall	f1-score	support
0	0.96	0.93	0.95	940
1	0.10	0.17	0.12	42
accuracy			0.90	982
macro avg	0.53	0.55	0.54	982
weighted avg	0.92	0.90	0.91	982

In [174]:



```
1 from sklearn.metrics import f1_score, precision_score, recall_score, confusion_matrix
2 # Print f1, precision, and recall scores
3 preScore = precision_score(y_test_labels,y_predicted_labels, average="macro")
4 recScore = recall_score(y_test_labels,y_predicted_labels, average="macro")
5 f1Score = f1_score(y_test_labels,y_predicted_labels, average="macro")
6 print(f"Precision Score = {preScore}")
7 print(f"Recall Score = {recScore}")
8 print(f"F1 Score = {f1Score}")
```

```
Precision Score = 0.530086114933288
```

```
Recall Score = 0.549290780141844
```

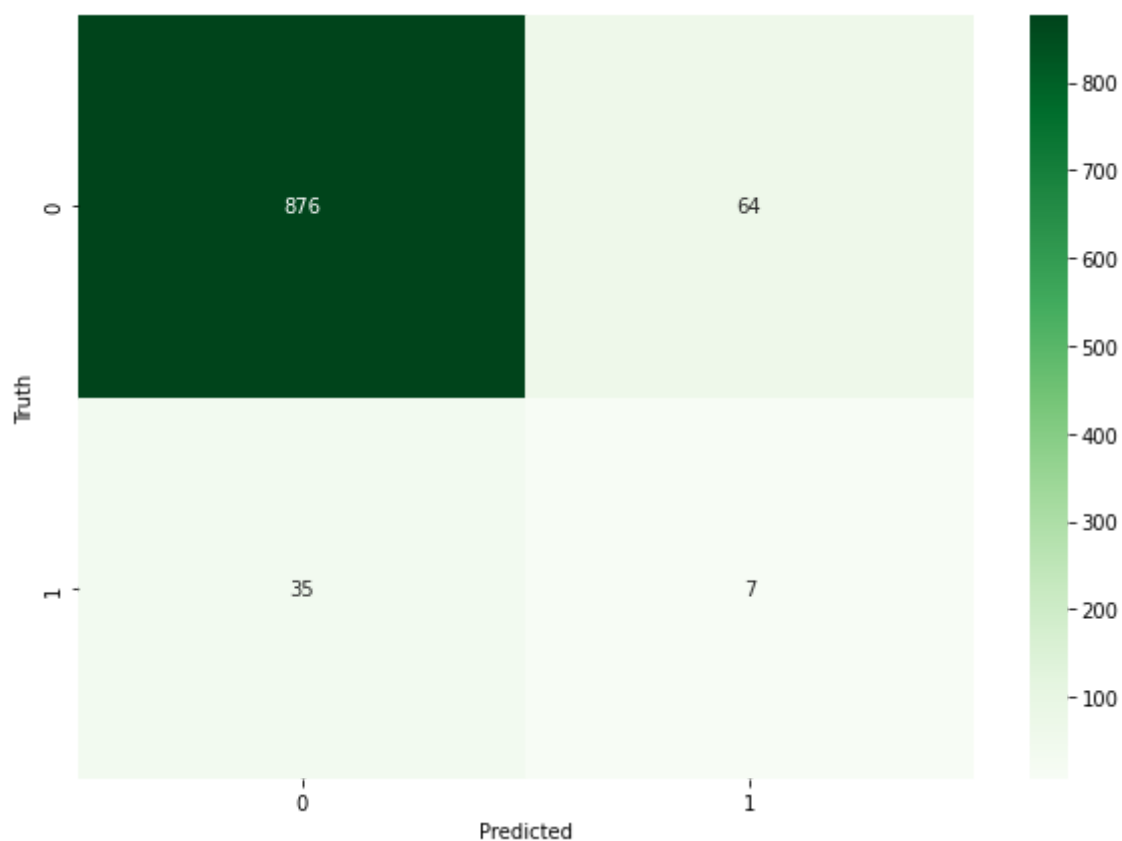
```
F1 Score = 0.5352046011961963
```

In [175]:

```
1 cm = tf.math.confusion_matrix(labels=y_test_labels,predictions=y_predicted_labels)
2 plt.figure(figsize=(10,7))
3 sns.heatmap(cm,annot=True,fmt='d',cmap='Greens')
4 plt.xlabel('Predicted')
5 plt.ylabel('Truth')
```

Out[175]:

Text(69.0, 0.5, 'Truth')



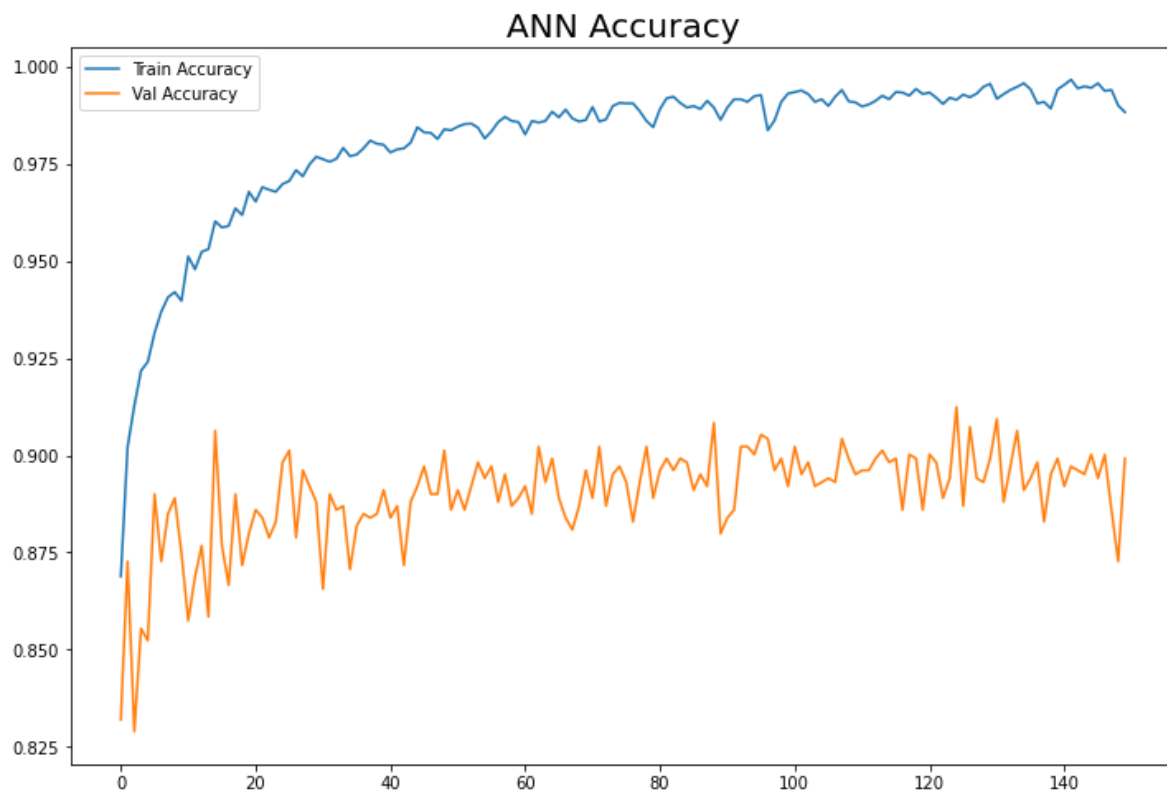
In [176]:



```

1  ## Plot History
2  fig = plt.figure(figsize=(12, 8))
3  plt.title('ANN Accuracy', size=20)
4  plt.plot(history.history['accuracy'], label="Train Accuracy")
5  plt.plot(history.history['val_accuracy'], label="Val Accuracy")
6  plt.legend()
7  plt.savefig("Visualizations\ANNAccuracyGraph.png")
8  plt.show()
9  model.evaluate(X_test,y_test)
10 score = np.round(model.evaluate(X_test, y_test, verbose=0)[1], 3)
11 print(f'Neural Network score      =====>>> {score}')
12 models_resultsANN = ({
13     'Model Name': 'Artificial Neural Network',
14     'Model Score': score,
15     'F1 Score': f1Score,
16     'Precision Score': preScore,
17     'Recall Score': recScore
18 })

```



```

31/31 [=====] - 0s 1ms/step - loss: 1.2418 - accuracy: 0.8992
Neural Network score      =====>>> 0.899

```

In [177]:



```
1 myalgorithms = {
2     'Logistic Regression': models_resultsLR,
3     'Gradient Boost': models_resultsGB,
4     'Decision Tree': models_resultsDT,
5     'Random Forest Classifier': models_resultsRC,
6     'K Nearest Neighbour': models_resultsKN,
7     'Artificial Neural Network': models_resultsANN
8 }
9 myalgorithms
```

Out[177]:

```
{'Logistic Regression': {'Model Name': 'Logistic Regression',
  'Model Score': 0.8360488798370672,
  'F1 Score': 0.15706806282722513,
  'Precision Score': 0.10067114093959731,
  'Recall Score': 0.35714285714285715},
 'Gradient Boost': {'Model Name': 'Gradient Boost',
  'Model Score': 0.9093686354378818,
  'F1 Score': 0.08247422680412371,
  'Precision Score': 0.07272727272727272,
  'Recall Score': 0.09523809523809523},
 'Decision Tree': {'Model Name': 'Decision Tree',
  'Model Score': 0.8686354378818737,
  'F1 Score': 0.1342281879194631,
  'Precision Score': 0.09345794392523364,
  'Recall Score': 0.23809523809523808},
 'Random Forest Classifier': {'Model Name': 'Random Forest',
  'Model Score': 0.9134419551934827,
  'F1 Score': 0.12371134020618557,
  'Precision Score': 0.10909090909090909,
  'Recall Score': 0.14285714285714285},
 'K Nearest Neighbour': {'Model Name': 'KNN',
  'Model Score': 0.869653767820774,
  'F1 Score': 0.15789473684210525,
  'Precision Score': 0.10909090909090909,
  'Recall Score': 0.2857142857142857},
 'Artificial Neural Network': {'Model Name': 'Artificial Neural Network',
  'Model Score': 0.899,
  'F1 Score': 0.5352046011961963,
  'Precision Score': 0.530086114933288,
  'Recall Score': 0.549290780141844}}}
```

In [178]:

```

1 myalgorithmsdf = pd.DataFrame(myalgorithms.values(),columns=['Model Name','Model Score']
2 myalgorithmsdf.sort_values(by='Model Score',ascending=False, inplace=True)
3 myalgorithmsdf

```

Out[178]:

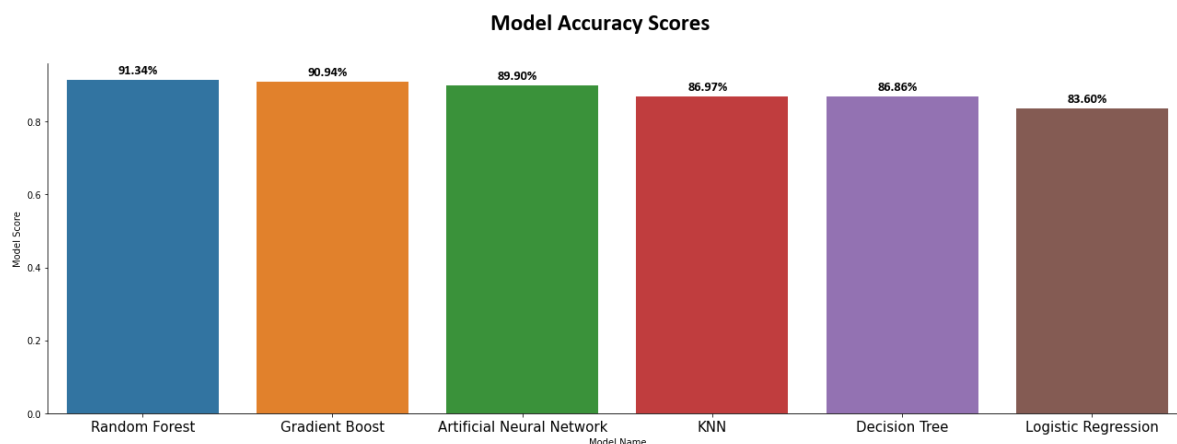
	Model Name	Model Score	F1 Score	Precision Score	Recall Score
3	Random Forest	0.913442	0.123711	0.109091	0.142857
1	Gradient Boost	0.909369	0.082474	0.072727	0.095238
5	Artificial Neural Network	0.899000	0.535205	0.530086	0.549291
4	KNN	0.869654	0.157895	0.109091	0.285714
2	Decision Tree	0.868635	0.134228	0.093458	0.238095
0	Logistic Regression	0.836049	0.157068	0.100671	0.357143

In [179]:

```

1 g = sns.catplot(x='Model Name', y='Model Score', data=myalgorithmsdf,
2               height=6, aspect=3, kind='bar', legend=True)
3 g.fig.suptitle('Model Accuracy Scores', size=25, y=1.1, **hfont)
4 ax = g.facet_axis(0,0)
5 ax.tick_params(axis='x', which='major', labelsize=15)
6 for p in ax.patches:
7     ax.text(p.get_x() + 0.27,
8           p.get_height() * 1.02,
9           '{0:.2f}%'.format(p.get_height()*100),
10          color='black',
11          rotation='horizontal',
12          size='x-large', **hfont)
13 plt.savefig("Visualizations\\ModelAccuracy.png")

```



In [180]:



```
1 # From the above diagram, it shows that Random Forest is the best model to use for the
```

SHAP

In [181]:



```
1 import shap
```

SHAP with RandomForestClassifier

In [183]:



```
1 X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=2)
2 sm = SMOTE(random_state=0)
3 X_train,y_train = sm.fit_resample(X_train,y_train)
4 scaler = StandardScaler()
5 X_train = scaler.fit_transform(X_train)
6 X_test = scaler.transform(X_test)
```

In [190]:



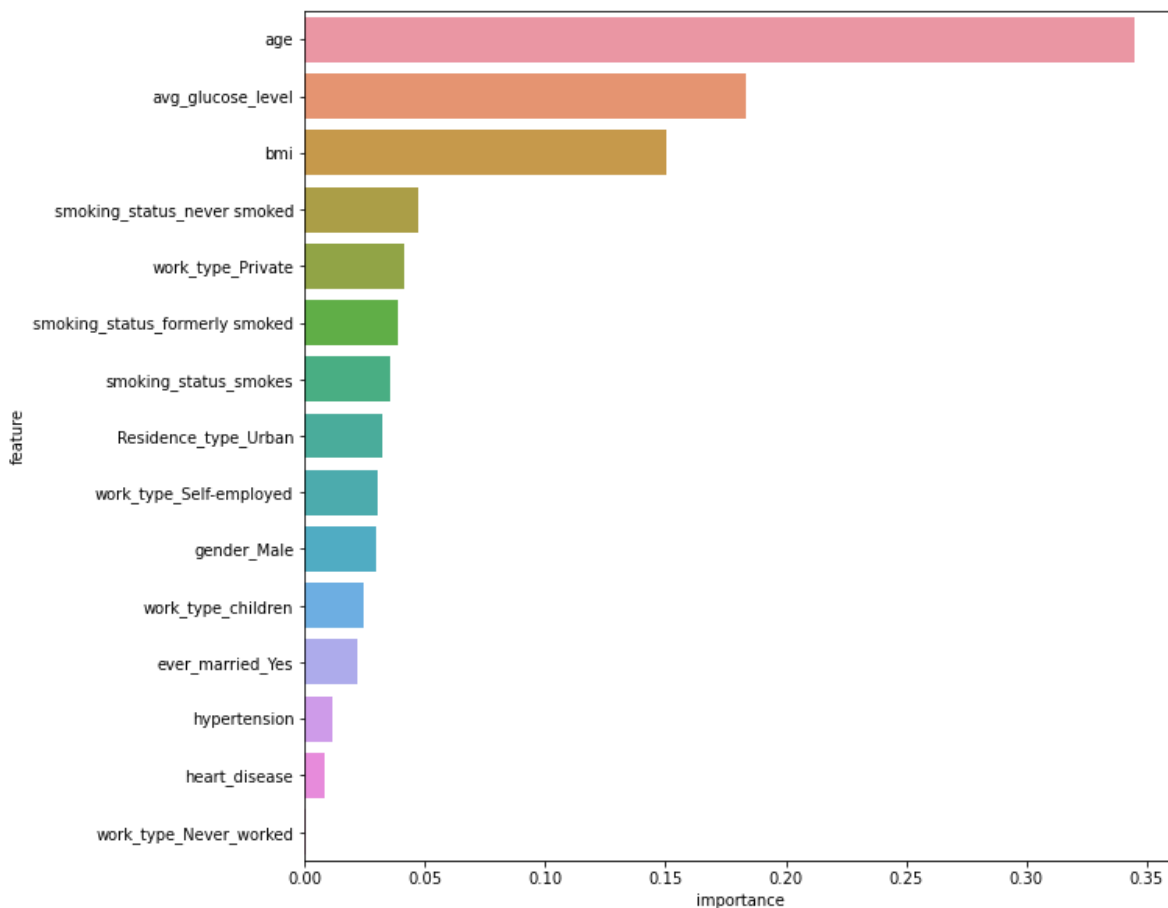
```
1 # Random Forest Feature Importance
```

In [185]:



```
1 RFModel = RandomForestClassifier(n_estimators=100,n_jobs=1)
2 RFModel.fit(X_train,y_train)
3 feature_importance_df = pd.DataFrame()
4 feature_importance_df['feature'] = X.columns
5 feature_importance_df['importance'] = RFModel.feature_importances_
6
7 feature_importance_df = feature_importance_df.sort_values('importance',ascending=False)
8 print('***Random Forest***')
9 plt.figure(figsize=(10,10))
10 sns.barplot(x='importance',y='feature',data=feature_importance_df[:15])
11 plt.savefig("Visualizations\ShapValueRandomForest.png")
12 plt.show()
```

Random Forest

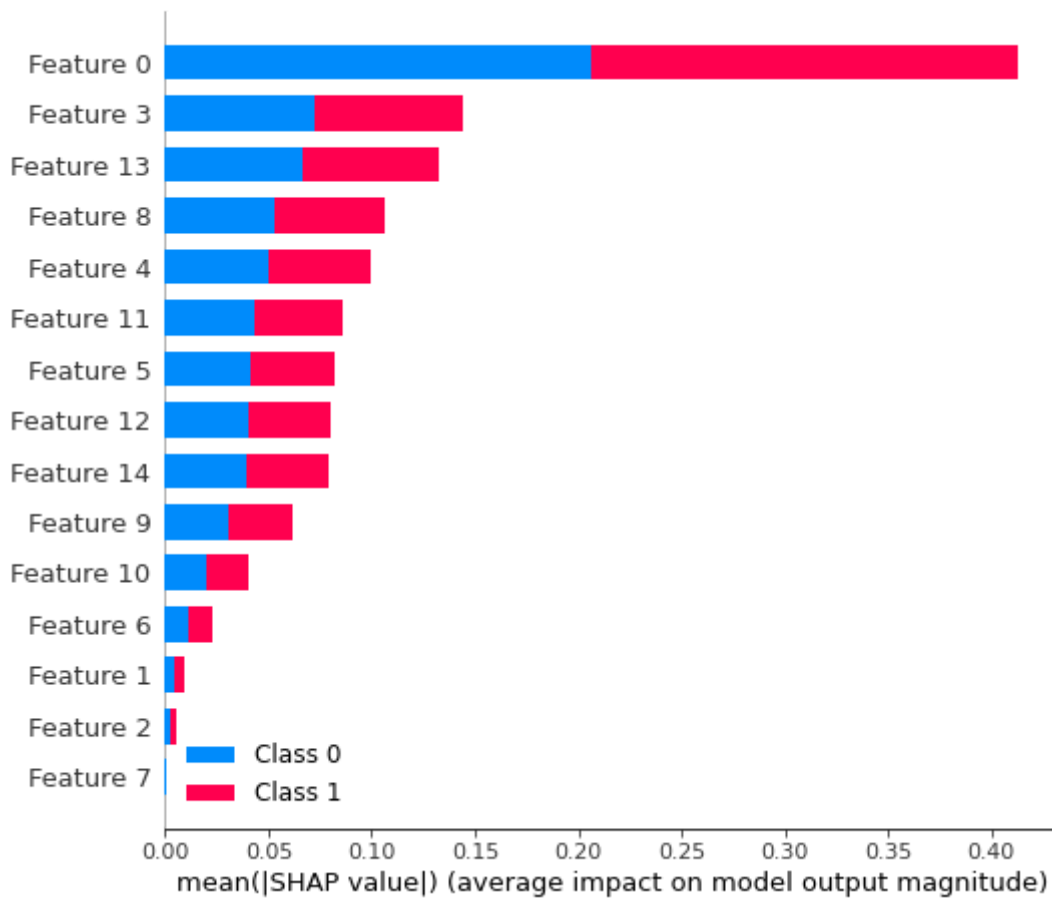


In [191]:



```
1 #Initialize SHAP Tree Explainer
2 explainer = shap.TreeExplainer(RFModel,model_output='margin')
3 shap_values = explainer.shap_values(X_test)
4
5 #Baseline Value
6 expected_value = explainer.expected_value
7 if isinstance(expected_value,list):
8     expected_value = expected_value[1]
9
10 print(f'Explainer Expected Value: {expected_value}')
11 shap.summary_plot(shap_values,X_test)
```

Explainer Expected Value: [0.50111998 0.49888002]



<Figure size 432x288 with 0 Axes>

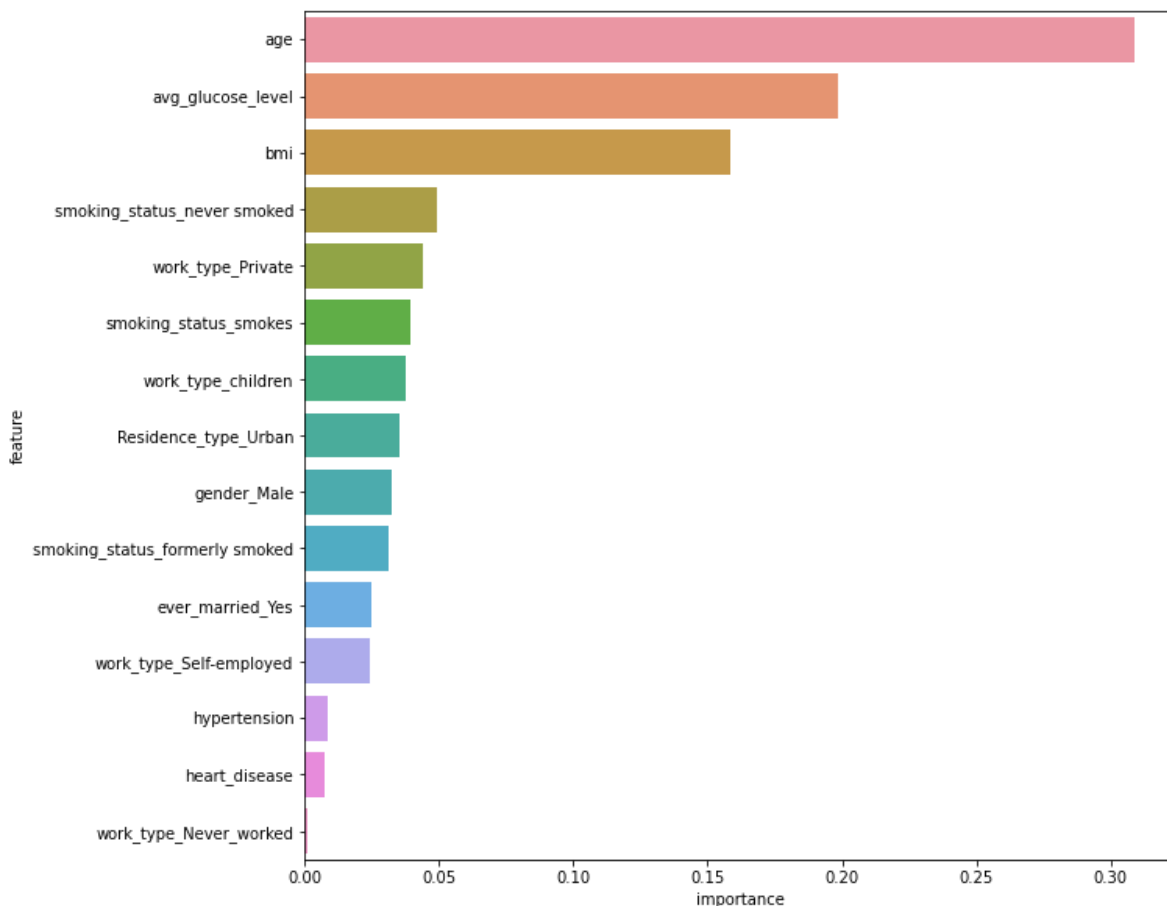
In []:

```
1 # Gradient Boost Feature Importance
```

In [193]:

```
1 GBModel = GradientBoostingClassifier(learning_rate=0.01,loss='exponential',max_depth=10)
2 GBModel.fit(X_train,y_train)
3 feature_importance_df = pd.DataFrame()
4 feature_importance_df['feature'] = X.columns
5 feature_importance_df['importance'] = GBModel.feature_importances_
6
7 feature_importance_df = feature_importance_df.sort_values('importance',ascending=False)
8 print('***Gradient Boost***')
9 plt.figure(figsize=(10,10))
10 sns.barplot(x='importance',y='feature',data=feature_importance_df[:15])
11 plt.savefig("Visualizations\ShapValueGradientBoost.png")
12 plt.show()
```

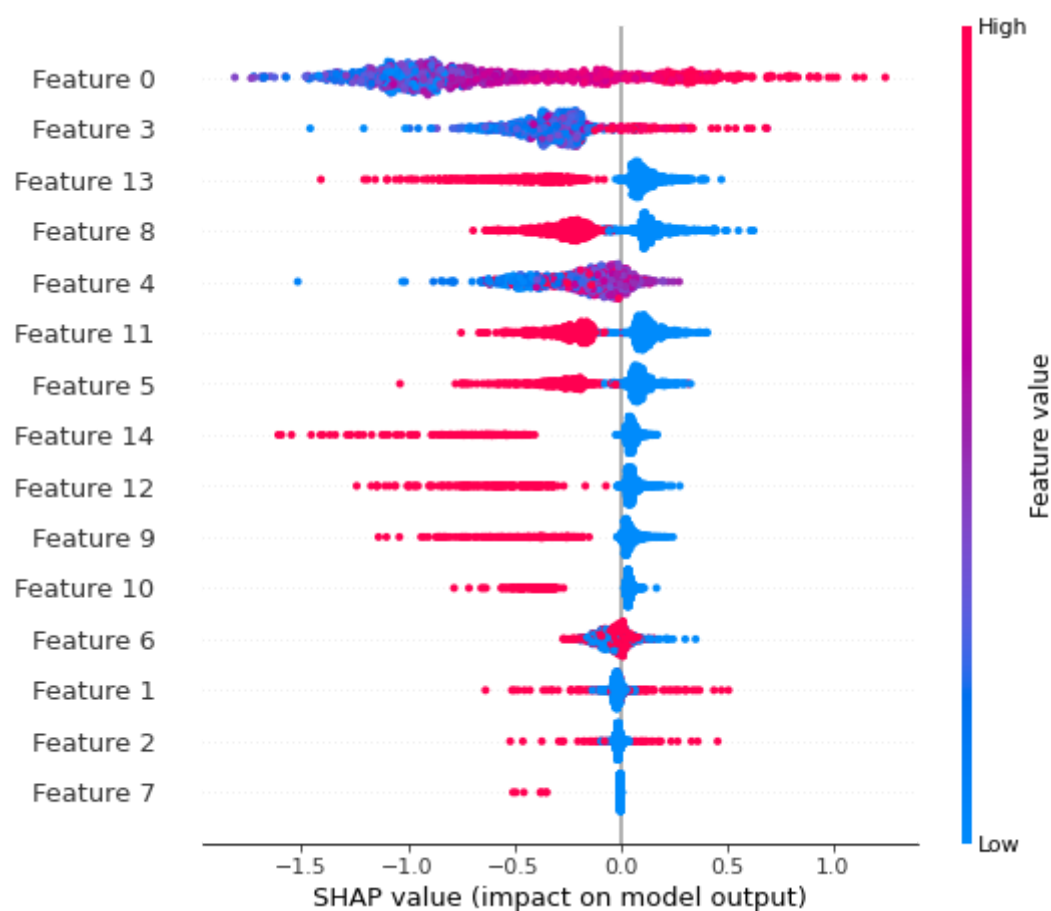
Gradient Boost



In [194]:

```
1 #Initialize SHAP Tree Explainer
2 explainer = shap.TreeExplainer(GBModel,model_output='margin')
3 shap_values = explainer.shap_values(X_test)
4
5 #Baseline Value
6 expected_value = explainer.expected_value
7 if isinstance(expected_value,list):
8     expected_value = expected_value[1]
9
10 print(f'Explainer Expected Value: {expected_value}')
11 shap.summary_plot(shap_values,X_test)
```

Explainer Expected Value: [8.28089334e-18]



<Figure size 432x288 with 0 Axes>

In []:



1	
---	--