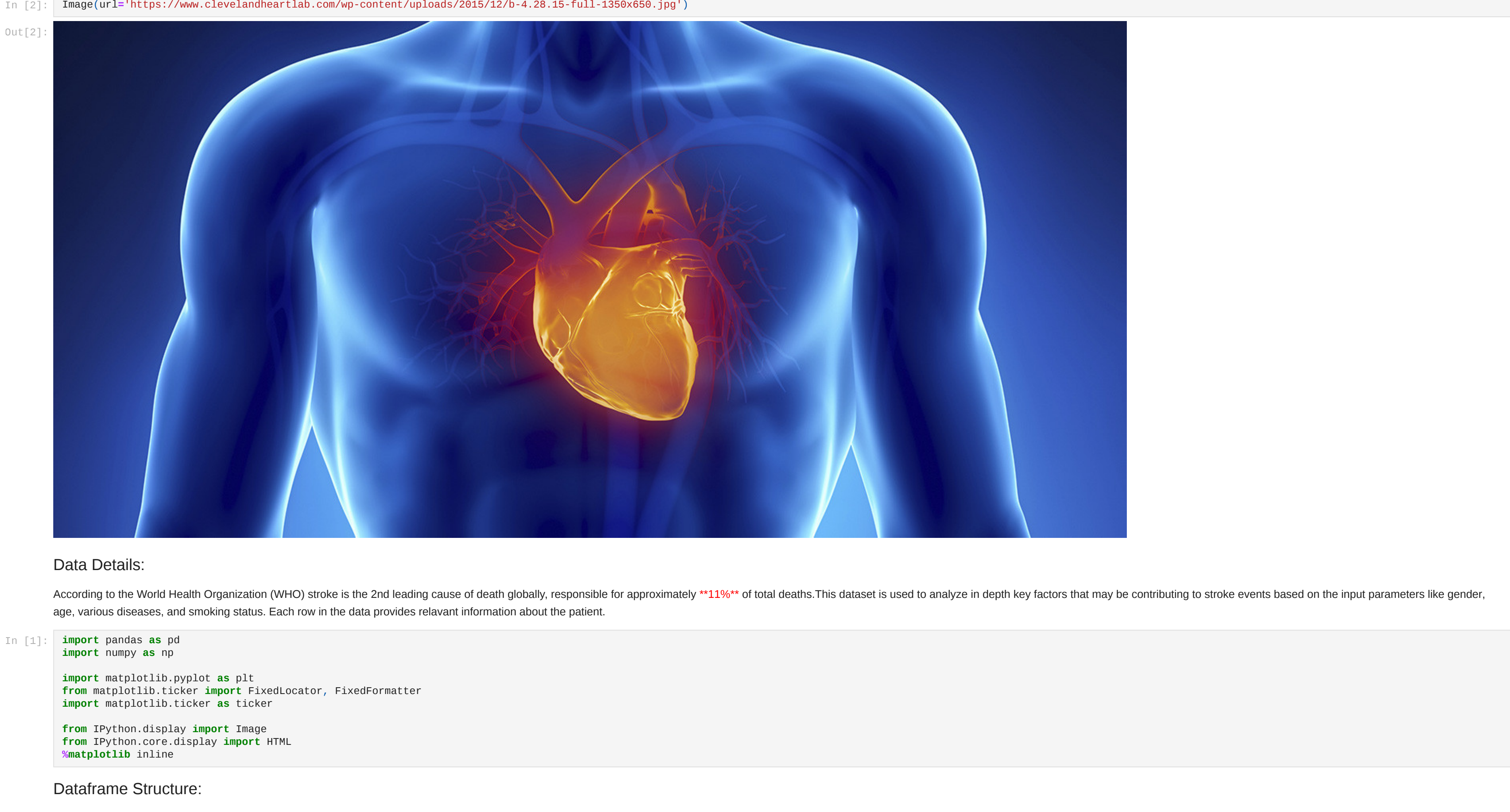


Stroke Event Analysis



Data Details:

According to the World Health Organization (WHO) stroke is the 2nd leading cause of death globally, responsible for approximately **~11%** of total deaths.This dataset is used to analyze in depth key factors that may be contributing to stroke events based on the input parameters like gender, age, various diseases, and smoking status. Each row in the data provides relevant information about the patient.

```
In [5]: import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
from matplotlib.ticker import FixedLocator, FixedFormatter
import matplotlib.ticker as ticker

from IPython.display import Image
from IPython.core.display import HTML
%matplotlib inline
```

Dataframe Structure:

This dataset contains over 5,000 datapoints regarding stroke events with attributes that are considered major factors in contributing to stroke instances around the globe. This healthcare dataset will be utilized along with data visualizations to demonstrate relationships or correlations between attributes. This analysis will demonstrate the impact of each attribute on stroke events, along with predictive trendspatterns.

Attributes:

1. id: unique identifier
 2. gender: "Male", "Female" or "Other"
 3. age: age of the patient
 4. hypertension: 0 if the patient doesn't have hypertension, 1 if the patient has hypertension
 5. heart_disease: 0 if the patient doesn't have any heart diseases, 1 if the patient has a heart disease
 6. ever_married: "No" or "Yes"
 7. work_type: "Children", "Govt_job", "Never_worked", "Private" or "Self-employed"
 8. Residence_type: "Rural" or "Urban"
 9. avg_glucose_level: average glucose level in blood
 10. bmi: body mass index
 11. smoking_status: "formerly smoked", "never smoked", "smokes" or "Unknown"
 12. stroke: 1 if the patient had a stroke or 0 if not
- ```
In [3]: # Loading dataframe with stroke healthcare data
stroke_data = pd.read_csv('healthcare-dataset-stroke-data.csv')
```

## Data Cleaning:

Original dataframe **stroke\_data** is refined/cleaned by removing ambiguous and unknown values. New dataframe is created **stroke\_refined** to create a more accurate analysis on stroke events. The affected attributes consist of **[smoking\_status], [bmi],** and **[gender]**. Dataframe query() function completes the desired filtering process using conditional statement, along with using dropna().

```
In [4]: stroke_refined = stroke_data.query('gender!="Other" & smoking_status!="Unknown"')
stroke_refined = stroke_refined.dropna()
```



## Hypertension, Heart Disease, and Strokes Associated With Gender:

This figure demonstrates the frequency of stroke events, hypertension, and heart disease attributes according to each gender group. A minor conclusion can be made based on this simple chart where women appear to have a higher stroke rate along with a higher hypertension rate. [Strokes Associated with Women](#). Not enough information is available on this chart to make solid conclusions but provides information regarding a possible relationship between strokes and hypertension. Heart disease is also a key factor and is presented with a higher frequency for men and could also be a main contributor to strokes. Hypertension is a attribute of interest and could be a main trigger for stroke events [The dangers of high blood pressure](#).



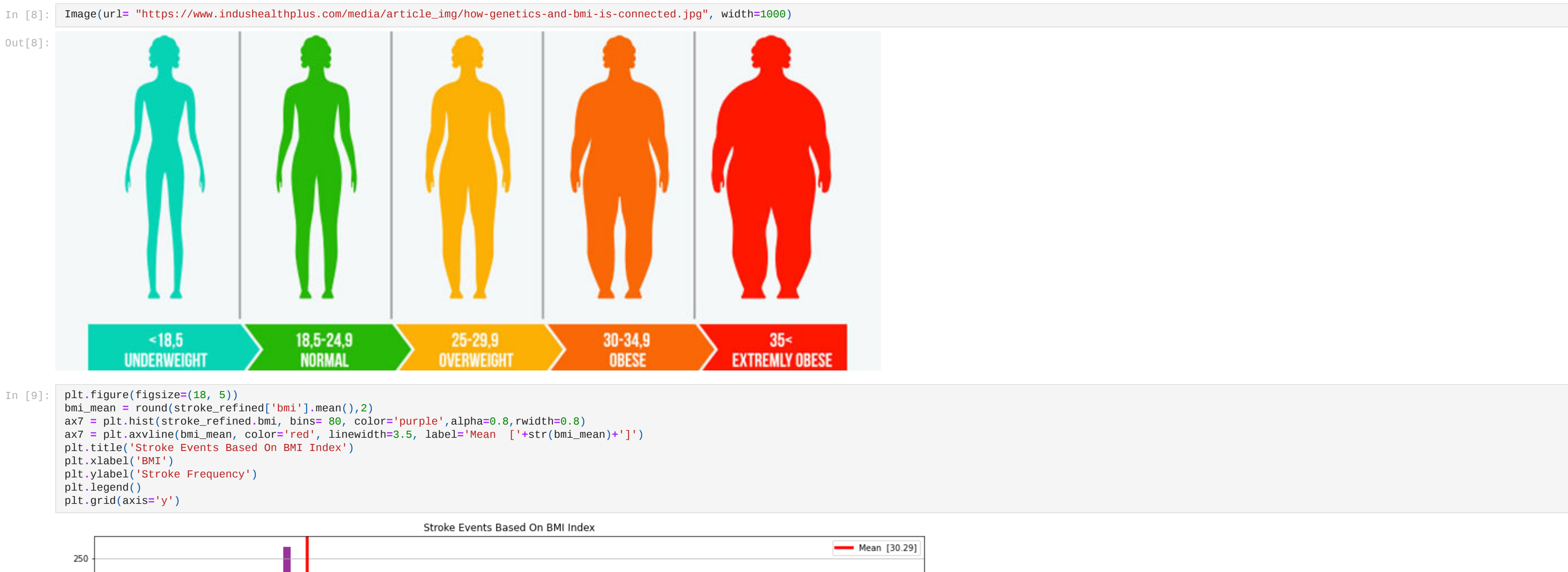
## Stroke Events Associated with Smoking Status:

Smoking status is observed and is known to be a major factor in contributing to many health conditions, including heart issues but findings have demonstrated little correlation between smokers and stroke events for this dataset. A smoking status bar chart is generated which provides frequencies for **hypertension, heart disease, and strokes** in connection to smoking status. The smoking status of **'smoker'** contains a small portion of the target data and doesn't seem to be a main contributor to stroke events. The smoking status of **'never smoked'** contains the largest frequencies in all target attributes with an obvious dominant frequency related to hypertension events.



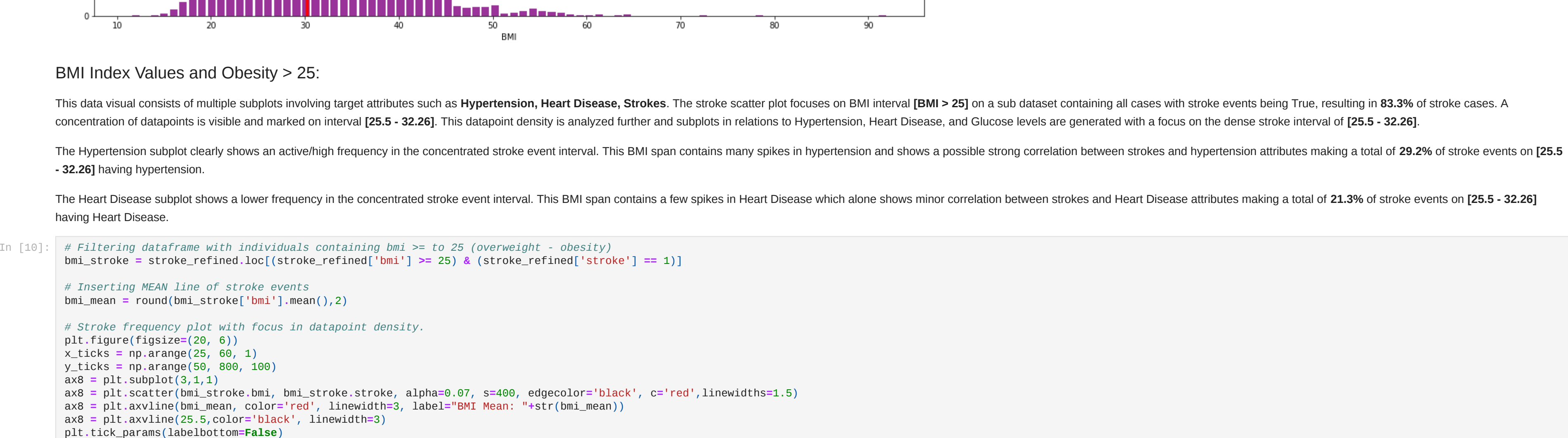
## Stroke Events Isolated By Smoker Status:

Based on the results from isolating smoker status and producing percentages per group, we can conclude that smoking isn't the major factor in stroke events in the dataset. Although this can be a factor in stroke events, smoking status alone doesn't provide enough evidence of being strongly correlated to strokes in this healthcare dataset. Here we can see 3 separate pie charts which show portions for **stroke vs non stroke events** in relation to smoking status and clearly demonstrates that **non stroke events** are dominant. Stroke events make a small percentage in every smoking status category and we can clearly see how most individuals with any type of smoking status don't experience a stroke.



## BMI Index Values and Obesity:

The BMI index dataset is analyzed with a histogram to demonstrate a distribution among all BMI values located in the healthcare dataset. Based on the BMI index charts, a conclusion has been made regarding a strong correlation between BMI values and strokes. The bmi attribute shows strong evidence of overweight - extreme obesity being a main contributor to stroke events **BMI value information**. The histogram generated from bmi values provides a clear view of bmi ranges in this dataset with a mean bmi of **30.29%** which clearly shows a large portion of individuals being obese and most being at the 27 BMI indicator. This histogram shows BMI distribution among entire dataset which includes BMI interval **[11.5 - 95]** which can be considered extreme values and will refine this interval further into this analysis.

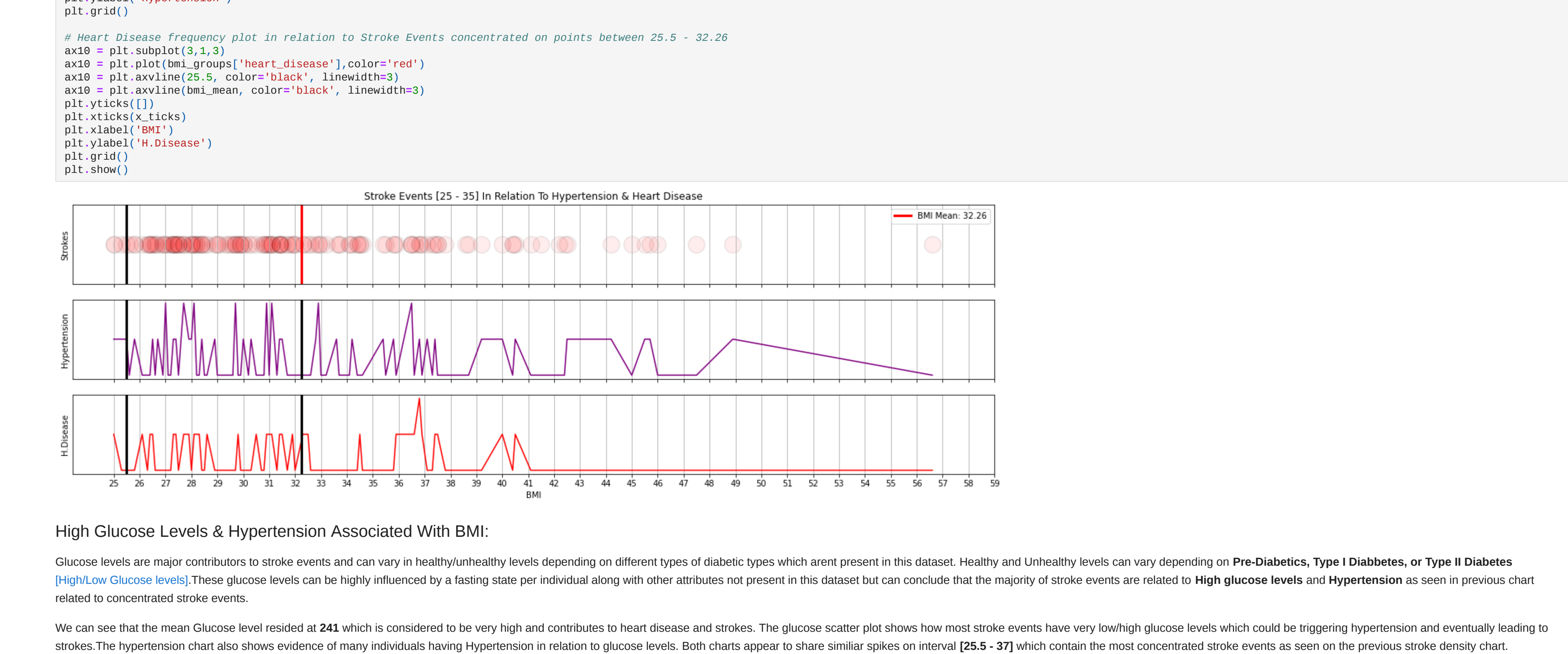


## BMI Index Values and Obesity > 25:

This data visual consists of multiple plots involving target attributes such as **Hypertension, Heart Disease, Strokes**. The stroke scatter plot focuses on BMI interval **[BMI > 25]** on a sub dataset containing all cases with stroke events being True, resulting in **83.3%** of stroke cases. A concentration of datapoints is visible and marked on interval **[28.5 - 32.26]**. This datapoint density is analyzed further and supports in relations to Hypertension, Heart Disease, and Glucose levels are generated with a focus on the dense stroke interval of **[25.5 - 32.26]**.

The hypertension scatter plot clearly shows an active/high frequency in the concentrated stroke event interval. This BMI span contains many spikes in hypertension and shows a possible strong correlation between strokes and hypertension attributes making a total of **29.2%** of stroke events on **[25.5 - 32.26]**.

The Heart Disease subplot shows a lower frequency in the concentrated stroke event interval. This BMI span contains a few spikes in Heart Disease which alone shows minor correlation between strokes and Heart Disease attributes making a total of **21.3%** of stroke events on **[25.5 - 32.26]** having Heart Disease.



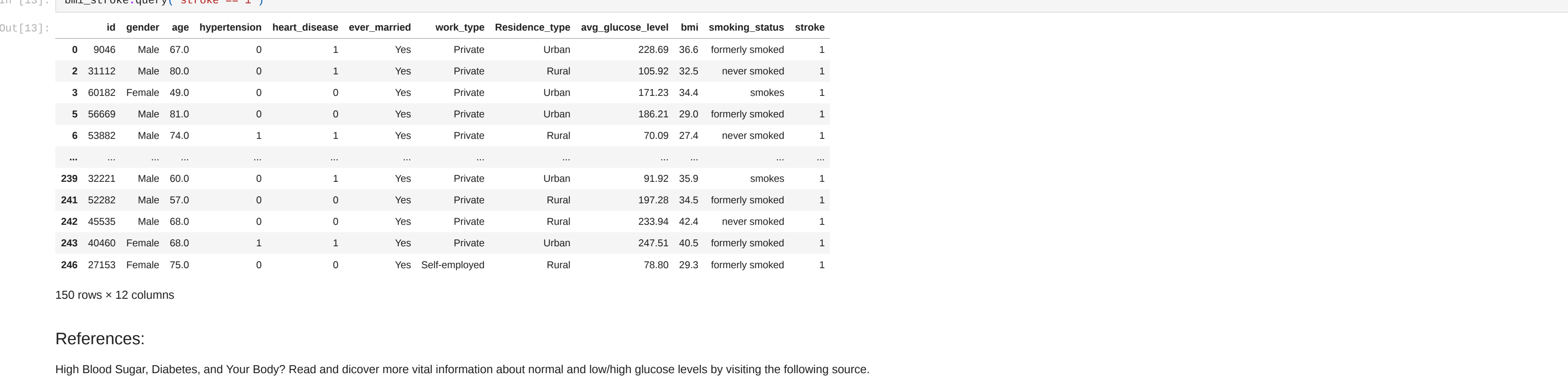
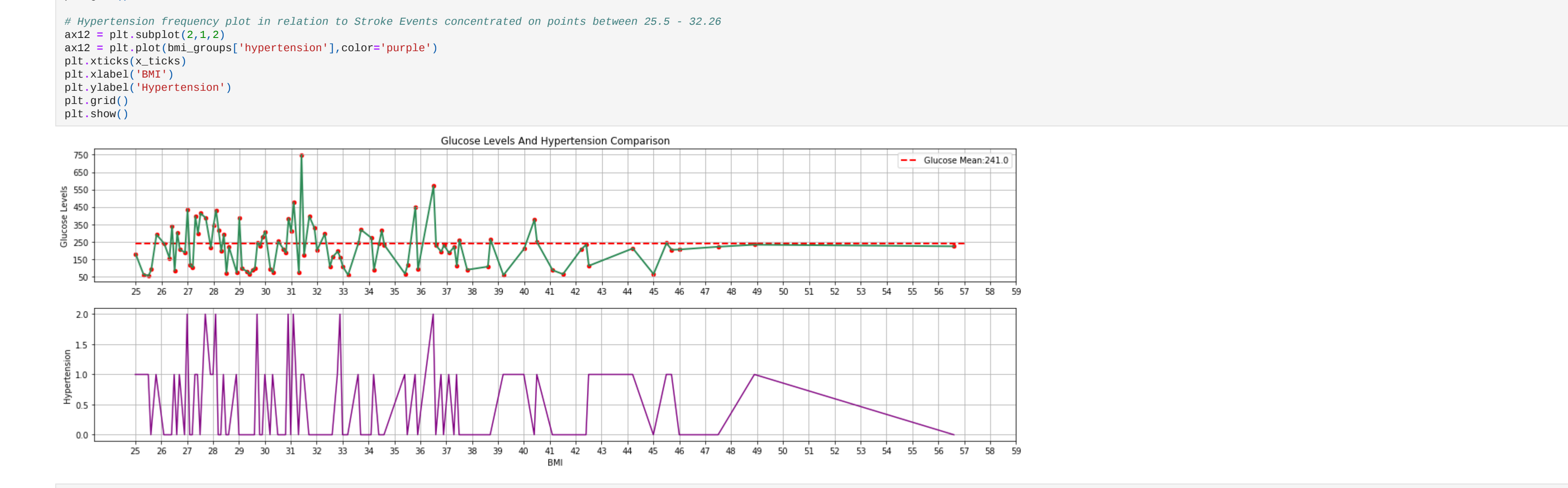
## High Glucose Levels & Hypertension Associated With BMI:

Glucose levels are major contributors to stroke events and can vary in healthy/unhealthy levels depending on different types of diabetic types which aren't present in this dataset. Healthy and unhealthy levels can vary depending on **Pre-Diabetics, Type I Diabetes, or Type II Diabetes** [High/Low Glucose Levels](#). These glucose levels can be highly influenced by a fasting state per individual along with other attributes not present in this dataset but can conclude that the majority of stroke events are related to **High glucose levels and Hypertension** as seen in previous chart related to concentrated stroke events.

We can see that the mean Glucose level resided at **241** which is considered to be very high and contributes to heart disease and strokes. The glucose scatter plot shows how most stroke events have very low/high glucose levels which could be triggering hypertension and eventually leading to strokes. The hypertension chart also shows evidence of many individuals having Hypertension in relation to glucose levels. Both charts appear to share similar spikes on interval **[25.5 - 37]** which contain the most concentrated stroke events as seen on the previous stroke density chart.

## Stroke Event Conclusion:

Hypertension and Glucose Level monitoring can be an effective way to prevent or predict stroke events but also have to take into consideration existing conditions such as **Heart Disease,Pre-Diabetics, Type I Diabetes, or Type II Diabetes**. These conditions are considered major influencers on healthy/unhealthy glucose levels which could lead to hypertension and eventually trigger strokes. We can conclude that hypertension is a key attribute in stroke events for this healthcare dataset but must take into consideration glucose levels along with pre-existing conditions.



## References:

High Blood Sugar, Diabetes, and Your Body? Read and discover more vital information about normal and low/high glucose levels by visiting the following source.

- Source: <https://www.webmd.com/diabetes/how-sugar-affects-diabetes>

What are normal glucose levels?

- Source: <https://www.vivianmason.org/whatarenormalbloodglucoselevels>

Learn more about the dangers regarding Hypertension/High blood pressure.

- Source: <https://www.cdc.gov/healthypeople/about/bmi/index.html>

BMI index values can be explained in depth by visiting the following sources.

- Source: <https://www.cdc.gov/healthypeople/about/bmi/index.html>

Women are at higher risk of having strokes? Visit the following sources which provide details behind why women might be at higher risk of experiencing strokes.

- Source 1: <https://www.stroke.org/en/about-stroke/stroke-risk-factors/women-have-a-higher-risk-of-stroke>

- Source 2: <https://utswmed.org/medblog/stroke-symptoms-women-risk/>