# 1 Organise Data

# 2 Data Visualization

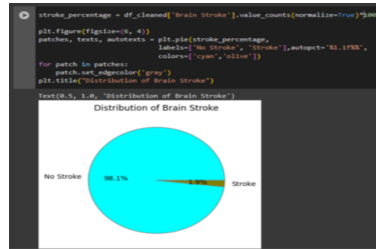1. How is the distribution of stroke observed across the dataset?



Figure 1: Pie Chart showing Distribution

In our dataset, stroke occurrence is highly prevalent, with 98.1% of records indicating stroke, while only 1.9% show no stroke. This significant imbalance necessitates careful consideration in predictive modelling to ensure accurate and reliable analysis.

2. How does the occurrence of stroke vary between genders within the dataset?



Figure 2: Vertical Bar Graph showing *Gender* with *Brain Stoke*

In our dataset, the occurrence of stroke varies between genders. Among males, 351 records indicate stroke, while 15,890 show no stroke. In comparison, 429 females have experienced stroke, with 23,777 records showing no stroke. This distribution suggests that females are slightly more affected by stroke in our dataset. Understanding such gender disparities is crucial for tailored preventive strategies and healthcare interventions. Further analysis can delve into the factors contributing to this discrepancy, potentially uncovering insights into gender-specific risk factors for stroke.

3. What is the total number of individuals who had experienced a stroke and also have hypertension?

Figure 3: Horizontal Bar graph showing *Hypertension* with *Brain Stroke*

The total number of individuals who experienced a stroke and have hypertension is approximately 5000. Notably, there's a higher frequency of stroke occurrence in individuals without hypertension, hinting at a potential inverse relationship between stroke and hypertension presence.

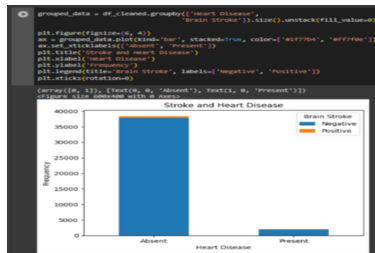4. Find the total number of people who beard the stroke have heart disease?



Figure 4: Vertical Bar graph showing *Heart Disease* with *Brain Stroke*

The data indicates fewer instances of heart disease among individuals who have experienced stroke compared to those without heart disease. Furthermore, individuals with stroke outnumber those without, suggesting a higher prevalence of stroke regardless of heart disease status.

5. How does the distribution of strokes differ between smokers non-smokers?



Figure 5: Vertical Bar graph showing *Smoking* with *Brain Stroke*

The distribution of strokes differs between smokers and non-smokers. The graph illustrates a higher count of individuals who never smoked, with a significant portion experiencing stroke. Conversely, fewer individuals who smoked are depicted, with a lower incidence of stroke. This suggests a potential association between non-smoking status and increased stroke occurrence.
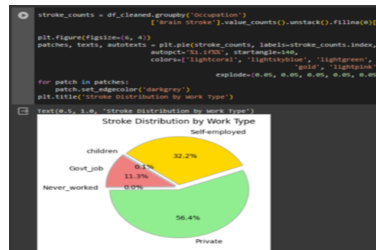
6. What is the occurrence of stoke based on occupation?



Figure 6: Pie Chart showing Distribution of *Brain Stroke* among *Occupation*

The graph displays the distribution of stroke occurrence across various occupation types. The majority (56.4%) are employed in the private sector, followed by self-employed individuals (32.2%). However, no records are present for individuals who have never worked, suggesting a potential data gap or exclusion from the dataset.

7. How does the occurrence of stroke vary based on the places individuals live?



Figure 7: Vertical Bar Graph showing *Brain Stroke* with *Residence Type*

The graph depicts a notable discrepancy in stroke occurrence based on individuals' places of residence. A higher count of individuals in both urban and rural areas is observed to have not experienced brain stroke, compared to those who have. This finding suggests a lower incidence of stroke among residents of both urban and rural areas, highlighting potential differences in lifestyle factors, healthcare access, or environmental influences between urban and rural populations.

8. What is the relationship between martial status and the occurrence of brain stroke?

Figure 8: Clustered Bar Graph showing *Brain Stroke* with *Martial Status*

The graph highlights a higher count of married individuals who have experienced stroke compared to those who have not. This suggests a potential association between marital status and increased risk of stroke occurrence among married individuals.

9. How does the level of glucose relate to the occurrence of stroke?



Figure 9: Vertical Bar Graph showing *Brain Stroke* with *Glucose Level*

The graph reveals a clear trend indicating a correlation between high levels of glucose and brain stroke occurrence. Individuals with elevated glucose levels are depicted as more likely to have experienced brain stroke, suggesting a potential association between high glucose levels and increased risk of stroke incidence. This finding underscores the importance of monitoring and managing glucose levels as part of stroke prevention strategies.

10. What is the association between body mass index (BMI) and the occurrence of stroke?

The graph illustrates a consistent prevalence of brain stroke across BMI values ranging from 20 to 45. This observation suggests a potential association between BMI within this range and the occurrence of brain stroke. Understanding this relationship could inform preventive measures targeting individuals within specific BMI ranges to mitigate stroke risk.

Figure 10: Strip Plot showing *Brain Stroke* with *BMI*

11. Is there any relation between Age and BMI?



(a)

(b)

(c)

(d)

Figure 11: K-Means

K-means clustering analysis was performed to explore the relationship between Age and BMI. The analysis revealed the presence of four distinct clusters, each representing a group of data points with similarities in Age and BMI values. By examining these clusters, we can gain insights into potential patterns or relationships between Age and BMI within our dataset

12. How does the occurrence of brain stroke vary with age and gender?

Analysis of the data reveals that female individuals show a higher susceptibility to brain stroke compared to males. However, when examining the relationship between age and brain stroke occurrence, no clear correlation emerges. This suggests that while gender may play a significant role in stroke risk, age alone may not be the sole determinant, indicating the influence of other factors in stroke occurrence among different demographic groups.
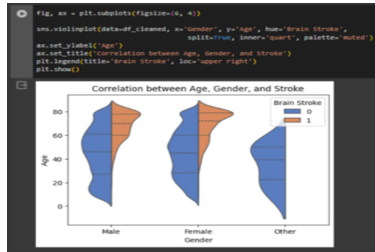
Figure 12: Violin Plot showing *Brain Stroke* with *Age* and *BMI*

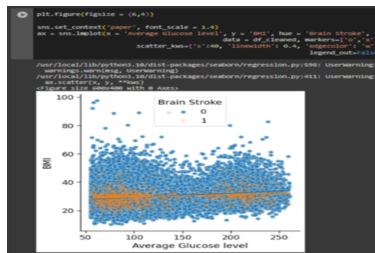13. How does the occurrence of stroke relate to both glucose level and body mass index (BMI)?



Figure 13: Scatter Plot showing *Glucose Level* with *BMI* based on *Brain Stroke*

The graph highlights that individuals with a BMI ranging from 25 to 35 and average glucose levels between 50 and 120 are particularly at risk of brain stroke. This pattern suggests a potential association between these factors and stroke occurrence. Understanding this relationship could inform targeted preventive strategies aimed at managing both glucose levels and BMI to reduce the risk of stroke.

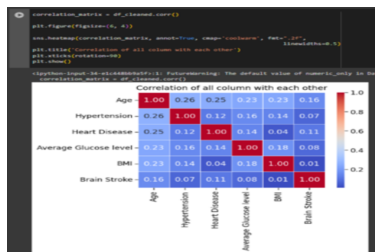14. How are all columns correlated with each other within the dataset?



Figure 14: HeatMap showing correlation among different variables

The analysis of correlations within the dataset reveals several notable associations. There is a significant correlation between brain stroke and age, indicating that age plays a crucial role in stroke occurrence. Additionally, a correlation is observed between brain stroke and heart disease, suggesting that individuals

with heart disease may be more predisposed to stroke. However, the correlation between BMI and brain stroke is weak, implying that while BMI may influence stroke risk, its impact is less pronounced compared to age and heart disease.

15. Will performing Principal Component Analysis (PCA) speed up the prediction process?
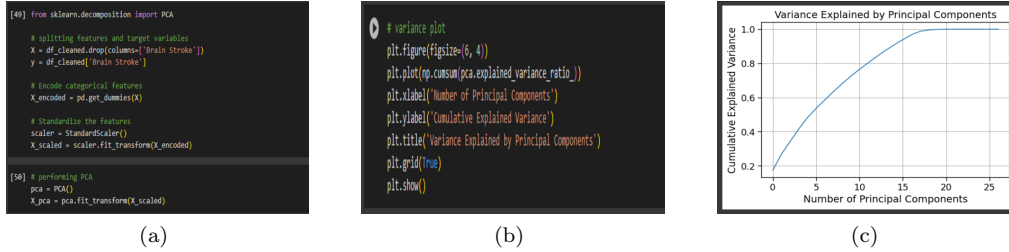


Figure 15: Principal Component Analysis

After performing Principal Component Analysis (PCA), it was observed that to achieve 80% variance explained, 11 principal components should be selected, as indicated by the variance plot. However, considering that the original dataset comprises only 10 features, the use of PCA in this context may not be particularly advantageous. In fact, selecting 11 principal components would result in an increase in the number of features rather than reducing dimensionality.
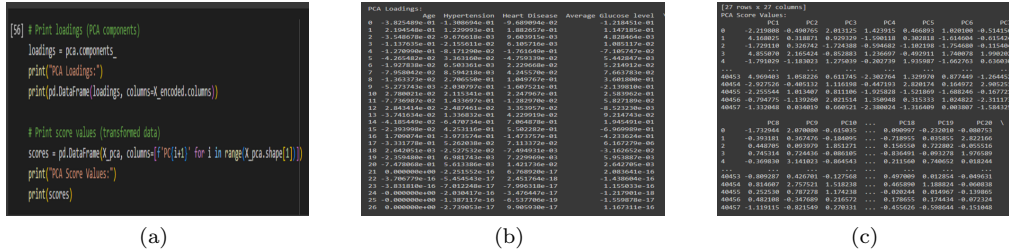


Figure 16: Principal Component Analysis

Therefore, while PCA is often utilized for dimensionality reduction and speeding up the prediction process, in this scenario, it may not provide any significant benefits and could potentially complicate the analysis without improving prediction speed.