

AIT 664 Part 2- Data Preparation & Information Modeling Report

G01394251

sai shivani rakam

Brain stroke prediction

In this part of data preparation and information modeling report, I will be collecting the data that is related to my search area and exploring it to clean the data, perform analysis and generate results.

1. Data Collection and importing.

The data for this analysis was collected from Kaggle. The dataset focuses on the key areas related to brain stroke, i.e., the main domains that cause brain strokes which is relevant to my research area.

The screenshot displays a Jupyter Notebook environment. At the top, the notebook is titled 'AIT-664' and 'Draft saved'. The menu bar includes 'File', 'Edit', 'View', 'Run', 'Add-ons', and 'Help'. On the right, there are buttons for 'Share', 'Save Version', and a version count of '0'. Below the menu, a toolbar contains icons for adding, deleting, copying, pasting, and running code. The main area shows a code cell with the following Python code:

```
[33]: bs = pd.read_csv('../input/full-filled-brain-stroke-dataset/full_data.csv')
bs
```

Below the code, a preview of the data is shown as a table with 11 columns: gender, age, hypertension, heart_disease, ever_married, work_type, Residence_type, avg_glucose_level, bmi, smoking_status, and stroke. The table displays rows 0 through 4980, with the last row being 4980. The 'stroke' column contains binary values (0 or 1). Below the table, it indicates '4981 rows x 11 columns'. There are buttons for '+ Code' and '+ Markdown'. Below the table, there is a code cell with the command 'bs.info()' and its output:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4981 entries, 0 to 4980
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   gender          4981 non-null   object
```

The screenshot shows a Kaggle notebook titled 'AIT-664' with a 'Draft saved' status. The interface includes a top bar with 'File', 'Edit', 'View', 'Run', 'Add-ons', and 'Help' menus, and a right bar with 'Share', 'Save Version', and a version count of '0'. A left sidebar contains icons for file management and search. The main area displays the output of the command `bs.info()`, which shows the dataset's structure: 4981 entries, 11 columns, and various data types. Below this, a code cell [36] contains Python code to print the unique values for several categorical columns. The output of this code is displayed below the cell, showing lists of unique values for 'gender', 'work_type', 'Residence_type', 'smoking_status', and 'ever_married'.

```
bs.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4981 entries, 0 to 4980
Data columns (total 11 columns):
#   Column              Non-Null Count  Dtype  
---  -
0   gender              4981 non-null  object  
1   age                 4981 non-null  float64  
2   hypertension         4981 non-null  int64  
3   heart_disease       4981 non-null  int64  
4   ever_married        4981 non-null  object  
5   work_type           4981 non-null  object  
6   Residence_type      4981 non-null  object  
7   avg_glucose_level   4981 non-null  float64  
8   bmi                 4981 non-null  float64  
9   smoking_status      4981 non-null  object  
10  stroke              4981 non-null  int64  
dtypes: float64(3), int64(3), object(5)
memory usage: 428.2+ KB

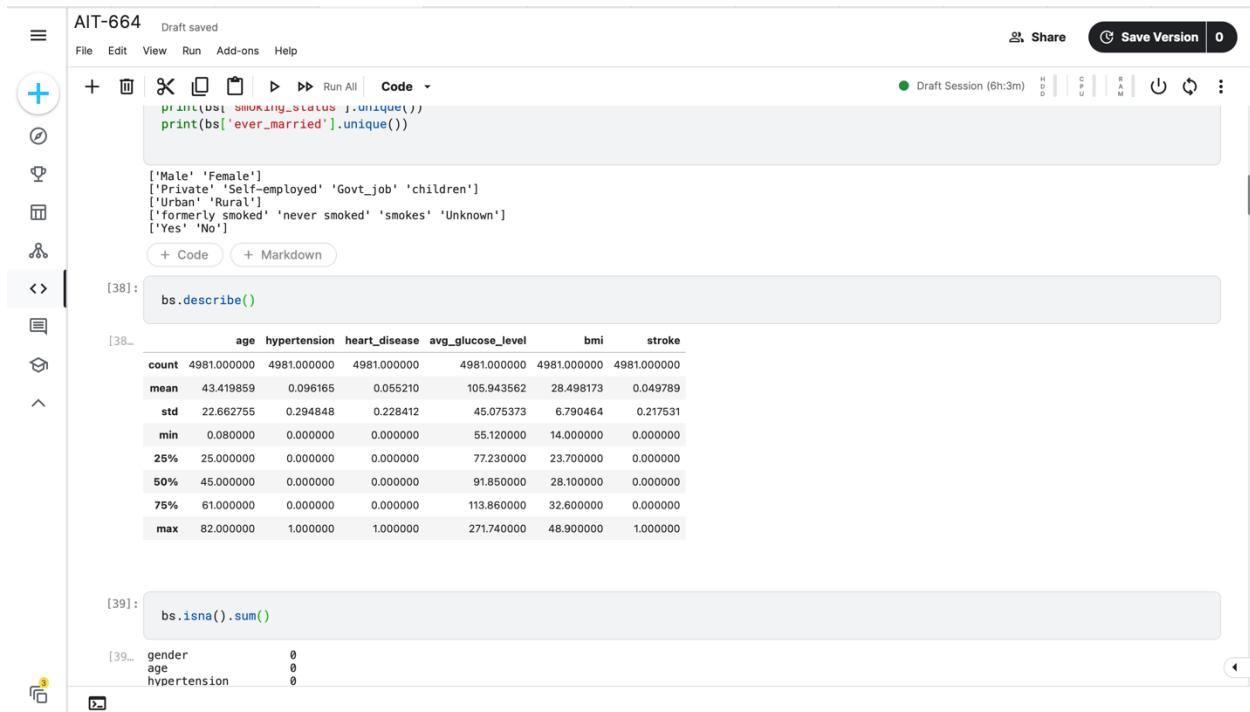
[36]: #seeing what are the categorical values
print(bs['gender'].unique())
print(bs['work_type'].unique())
print(bs['Residence_type'].unique())
print(bs['smoking_status'].unique())
print(bs['ever_married'].unique())

['Male' 'Female']
['Private' 'Self-employed' 'Govt_job' 'children']
['Urban' 'Rural']
['formerly smoked' 'never smoked' 'smokes' 'Unknown']
['Yes' 'No']
```

I used Kaggle notebook to explore and run required python code for my analysis. I imported the dataset with all the libraries. We can see that the dataset contains 4981 records and 11 features. The dataset has 11 columns. The above output tells the column name and the datatype of the column and the categorical values. The columns in this dataset are gender, age, hypertension, heart_disease, ever_married, work_type, residence_type, avg_glucose_level, bmi, smoking_status and stroke.

2. Data Preparation/cleaning

I checked for missing data and duplicates, there is no missing or duplicate data in the dataset. There is a categorical data in the smoke column that says UNKNOWN, I had a thought of changing it to never_smoked status but it could be misleading the original data. I also generated the summary statistics for the numerical data.



The screenshot shows a Jupyter Notebook interface with the following components:

- Top Bar:** "AIT-664" and "Draft saved". On the right, there are "Share", "Save Version", and a counter "0".
- Menu Bar:** File, Edit, View, Run, Add-ons, Help.
- Left Sidebar:** Contains icons for file operations (plus, minus, copy, paste, run, code, markdown) and a search icon.
- Code Editor:**
 - Cell [37]:

```
print(bs['smoking_status'].unique())
print(bs['ever_married'].unique())
```
 - Cell [38]:

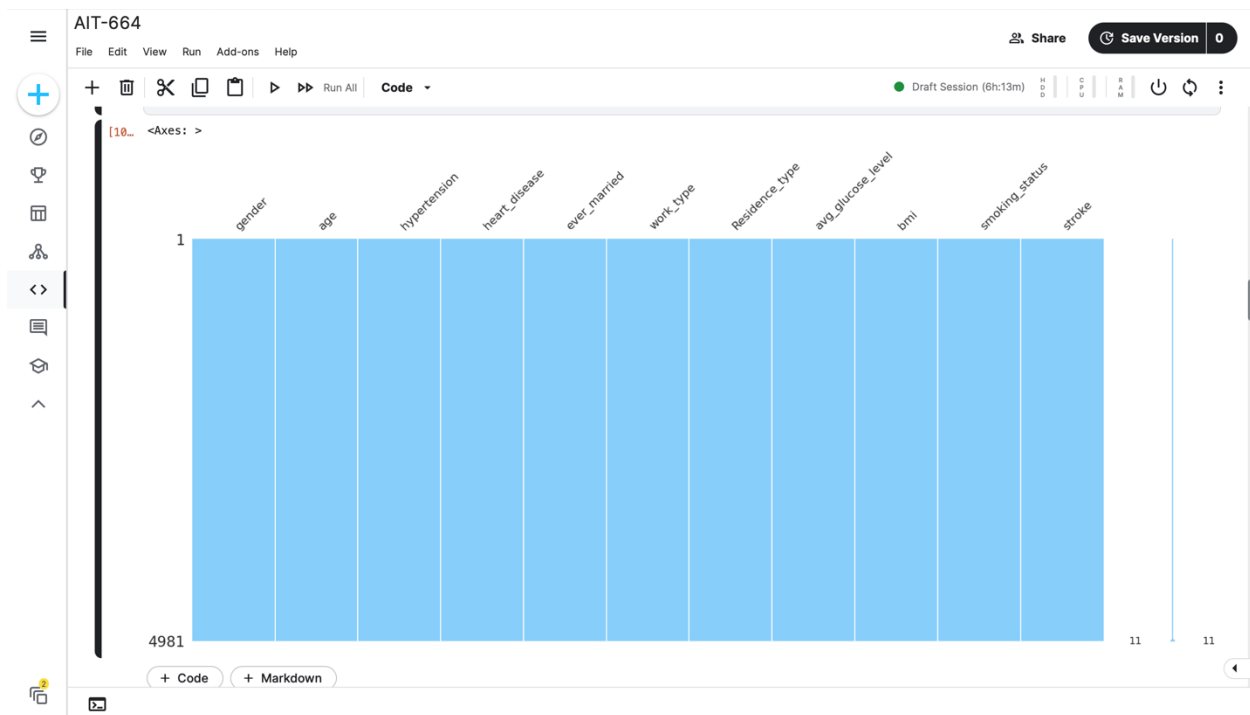
```
bs.describe()
```
 - Cell [39]:

```
bs.isna().sum()
```
- Output for Cell [38]:** A summary statistics table for the 'bs' dataset.

	age	hypertension	heart_disease	avg_glucose_level	bmi	stroke
count	4981.000000	4981.000000	4981.000000	4981.000000	4981.000000	4981.000000
mean	43.419859	0.096165	0.055210	105.943562	28.498173	0.049789
std	22.662755	0.294848	0.228412	45.075373	6.790464	0.217531
min	0.080000	0.000000	0.000000	55.120000	14.000000	0.000000
25%	25.000000	0.000000	0.000000	77.230000	23.700000	0.000000
50%	45.000000	0.000000	0.000000	91.850000	28.100000	0.000000
75%	61.000000	0.000000	0.000000	113.860000	32.600000	0.000000
max	82.000000	1.000000	1.000000	271.740000	48.900000	1.000000

Output for Cell [39]:

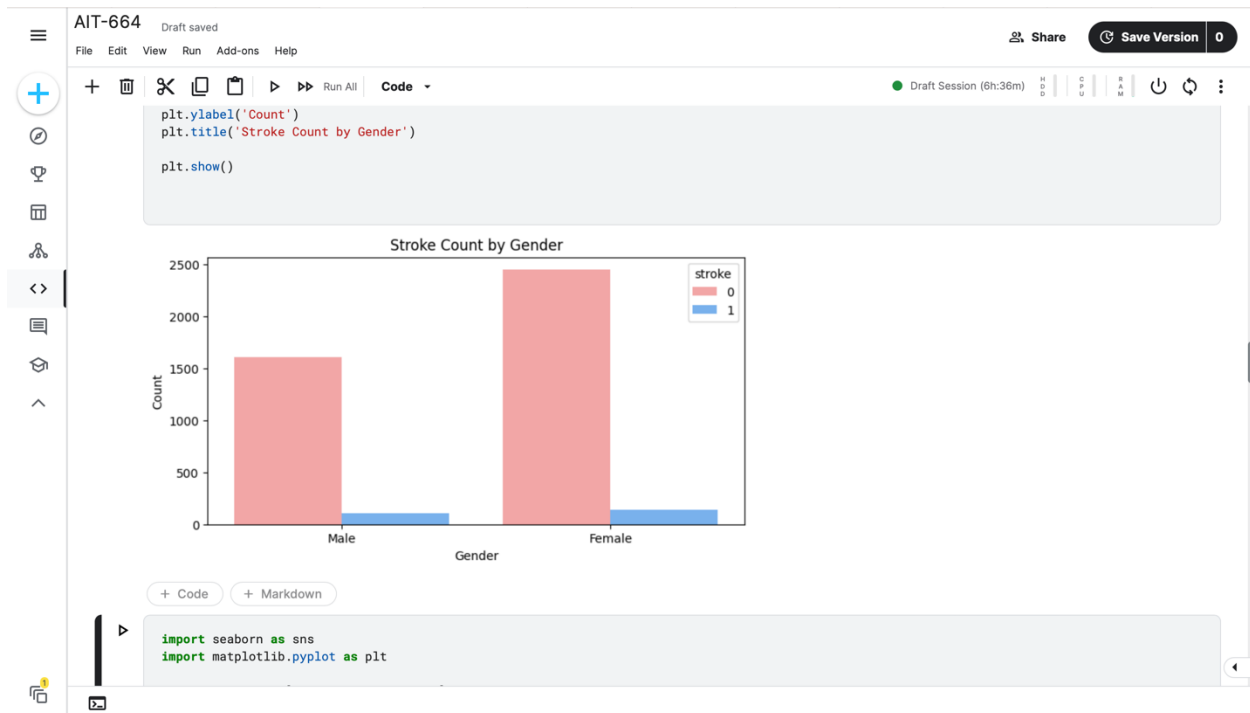
gender		0
age		0
hypertension		0

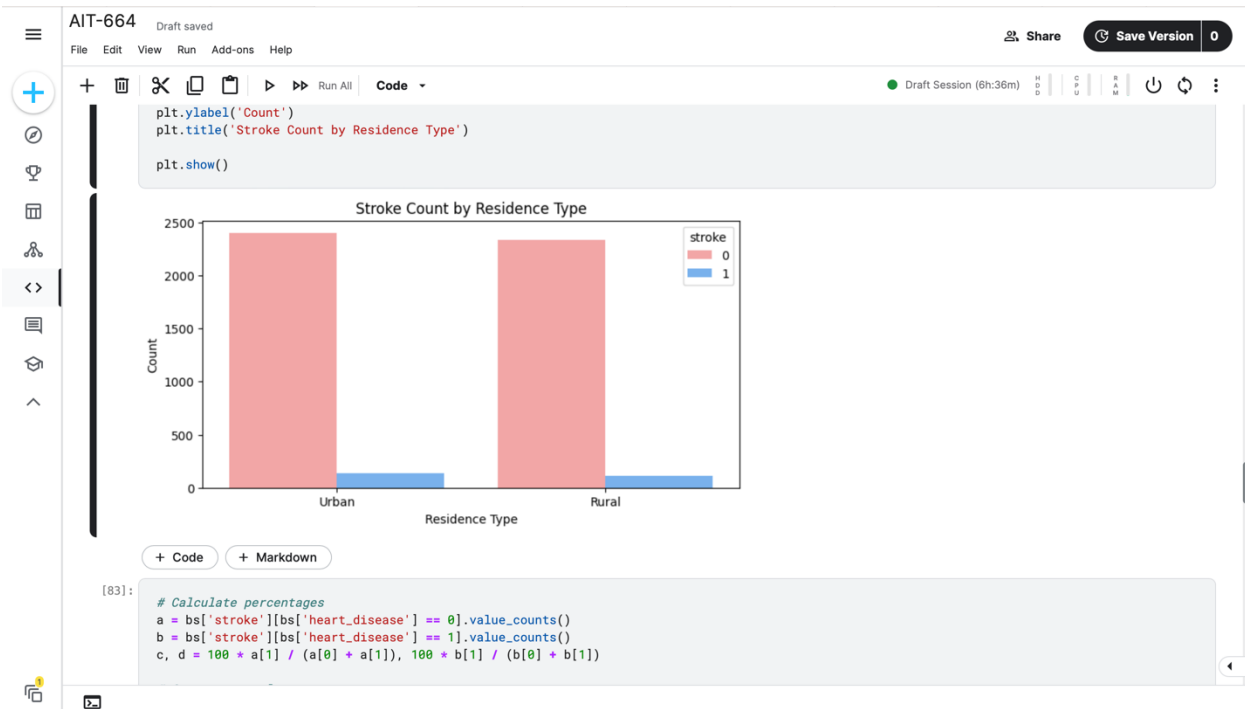
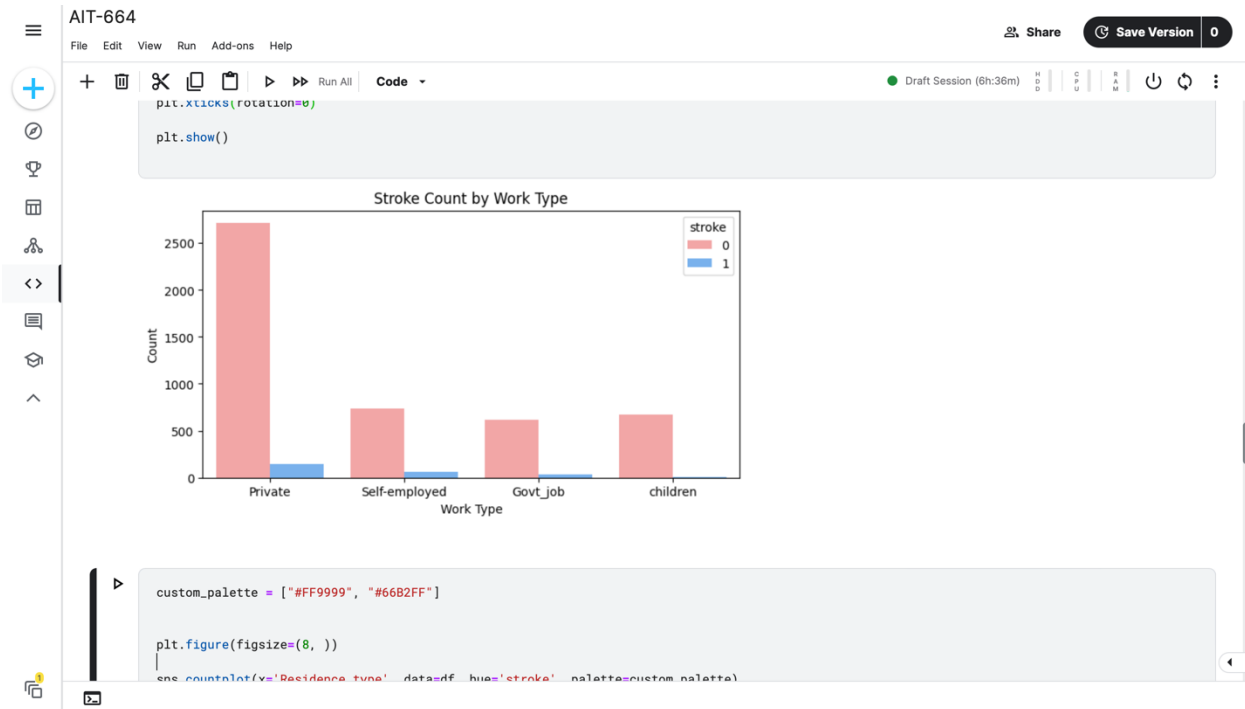


3. Exploratory Analysis

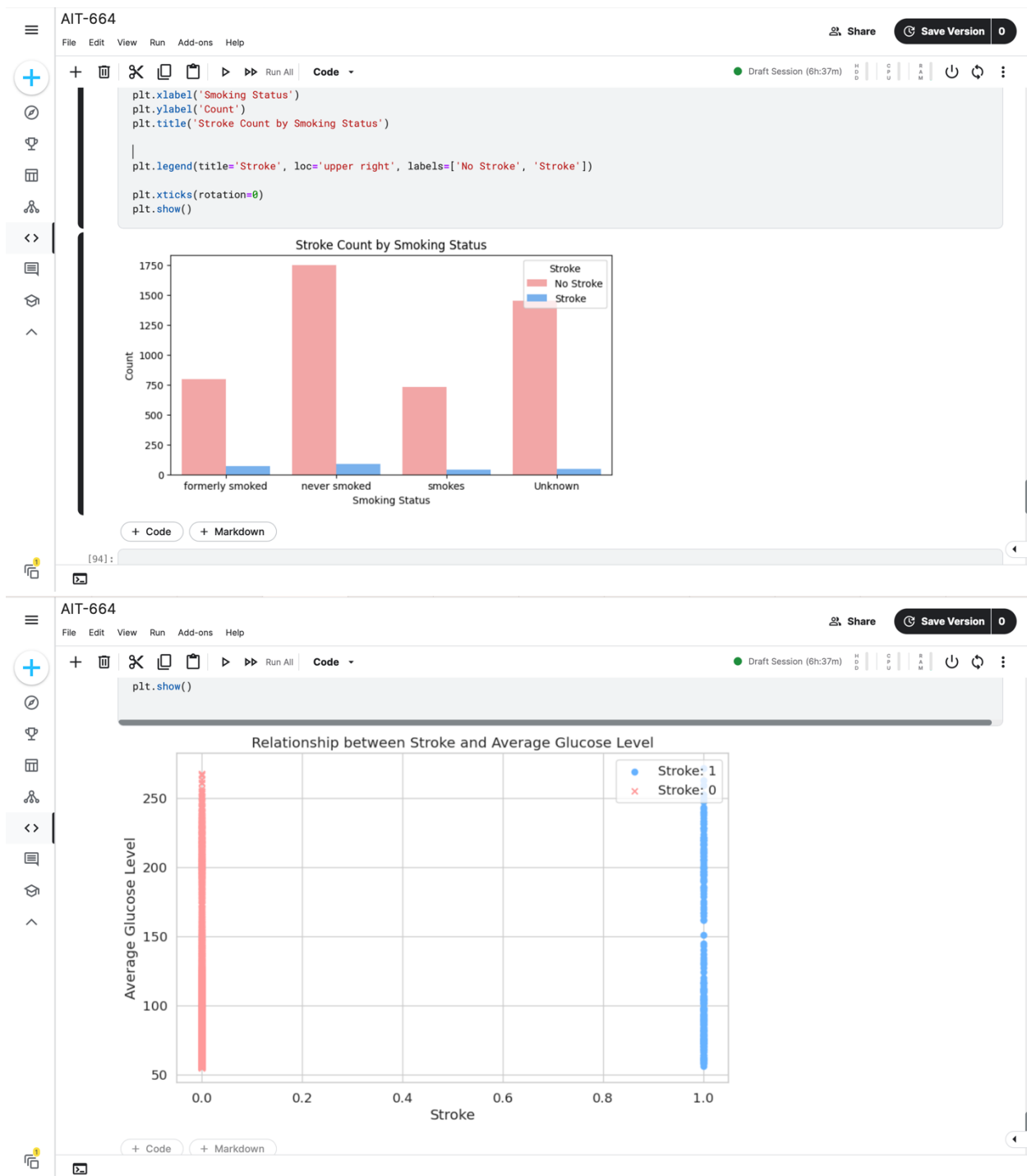
I performed the EDA correlating the features of the dataset with the target variable(stroke) this gave me surprising results.











Observation from outputs:

- It is observed that older people are more prone to stroke.
- Females are at risk of getting stroke than male.
- Strokes are also more likely to occur to people who have high glucose level compared to the ones with normal glucose level.

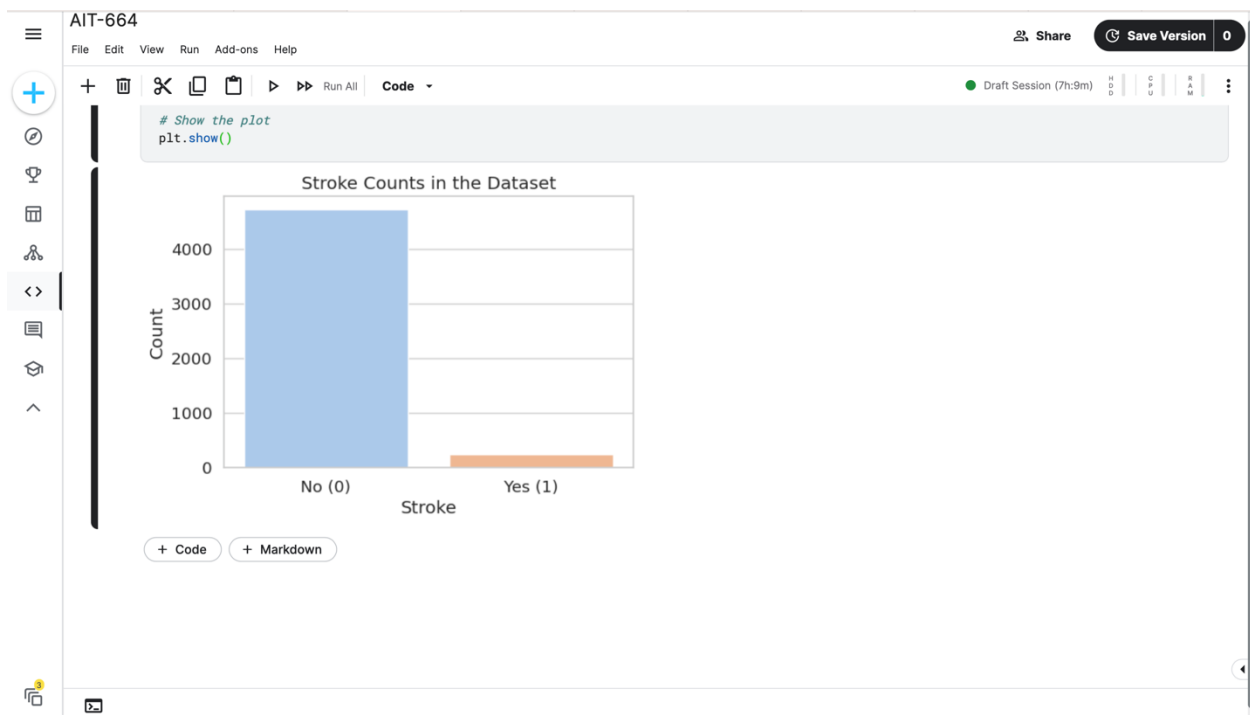
- Marriage also seems to be responsible for increasing the risk for stroke.
- there are also more percentage of people getting a stroke if they have hypertension and are suffering from any kind of heart disease.
- Smokers have higher risk of having a Brain Stroke.
- Urban people have a little more risk of stroke than rural.

4. Preprocessing/Modeling

In this step of preprocessing and modelling the data I further analyzed the data. I observed that the data is not balanced. Modelling the data that is not balanced would be challenging, it is a good decision to balance the data.

To balance it we can either be done by down sampling or over sampling, here I Performed down sampling on the majority class by randomly selecting 5% of the data. This down sampling technique is commonly used in machine learning to address class imbalance issues, where one class is significantly underrepresented compared to the other class. By creating a balanced subset, it can improve the performance of algorithms, especially in classification.

I utilized a Random Forest Classifier for my predictions, starting with a basic model and then refining it through hyperparameter tuning using RandomizedSearchCV. The dataset was split into training and testing sets for model evaluation.



AIT-664 Failed to save draft.

File Edit View Run Add-ons Help

Share Save Version 0

Draft Session (8h:40m)

```
# Now you have clf_gini initialized with max_depth=5 and random_state=0
```

[21] DecisionTreeClassifier

```
DecisionTreeClassifier(max_depth=5, random_state=0)
```

y_pred_rs = clf_gini.predict(X_test_rs)

```
print(confusion_matrix(y_test_rs, y_pred_rs))
print('The report is:\n{}'.format(classification_report(y_test_rs, y_pred_rs)))
```

[[59 24]
[18 63]]

The report is:

	precision	recall	f1-score	support
0	0.77	0.71	0.74	83
1	0.72	0.78	0.75	81
accuracy			0.74	164
macro avg	0.75	0.74	0.74	164
weighted avg	0.75	0.74	0.74	164

+ Code + Markdown

[214]:

```
rfc = RandomForestClassifier()
rfc.fit(X_train_rs, y_train_rs)
```

[21] RandomForestClassifier

```
RandomForestClassifier()
```

AIT-664 Failed to save draft.

File Edit View Run Add-ons Help

Share Save Version 0

Draft Session (8h:43m)

```
estimator: RandomForestClassifier
RandomForestClassifier()
RandomForestClassifier()
```

y_pred_rfc_random = rfc_random.predict(X_test_rs)

```
print(confusion_matrix(y_test_rs, y_pred_rfc_random))
print('The accuracy is: {:.4f}'.format(accuracy_score(y_test_rs, y_pred_rfc_random)))
print('The report is:\n{}'.format(classification_report(y_test_rs, y_pred_rfc_random)))
```

[[57 26]
[10 71]]

The accuracy is: 0.7885

The report is:

	precision	recall	f1-score	support
0	0.85	0.69	0.76	83
1	0.73	0.88	0.80	81
accuracy			0.78	164
macro avg	0.79	0.78	0.78	164
weighted avg	0.79	0.78	0.78	164

+ Code + Markdown

Observation:

- Random Under-Sampling: After balancing the dataset, both classes (0 and 1) had 248 instances each.
- Basic Random Forest Model: The initial Random Forest model achieved an accuracy of 74% on the test set, with a precision of 72% for stroke patients and 77% for non-stroke

patients. Recall was 78% for stroke patients and 71% for non-stroke patients. F1-scores were 0.75 for stroke patients and 0.74 for non-stroke patients.

- **Tuned Random Forest Model:** The hyperparameter-tuned Random Forest model showed significant improvement. It achieved an accuracy of 78%, with a precision of 73% for stroke patients and 85% for non-stroke patients. Recall was 88% for stroke patients and 69% for non-stroke patients. F1-scores were 0.80 for stroke patients and 0.76 for non-stroke patients.

The Random Forest model, especially after hyperparameter tuning, demonstrated promising results in predicting strokes. With an accuracy of 78%, the model showed strong precision and recall values, indicating its ability to correctly identify stroke patients while minimizing misclassifications. Further fine-tuning and exploration of different algorithms could potentially enhance the model's performance, making it a valuable tool in identifying patients at risk of strokes.

Conclusion:

In this data preparation and information modeling report, a comprehensive analysis was conducted on a dataset sourced from Kaggle focusing on brain stroke-related domains. The dataset was meticulously processed, normalized, and cleaned to ensure its quality and reliability for analysis. Notable observations emerged, revealing crucial insights into stroke risk factors. Key findings indicated that older age, female gender, high glucose levels, marriage, hypertension, heart diseases, and smoking are all significant contributors to stroke risk. An initial Random Forest model achieved a respectable accuracy of 74%. However, significant enhancements were achieved through hyperparameter tuning.

Reference:

Akbasli, I. T. (2022, July 16). *Brain stroke prediction dataset*. Kaggle.
<https://www.kaggle.com/datasets/zzettrkalpakbal/full-filled-brain-stroke-dataset>