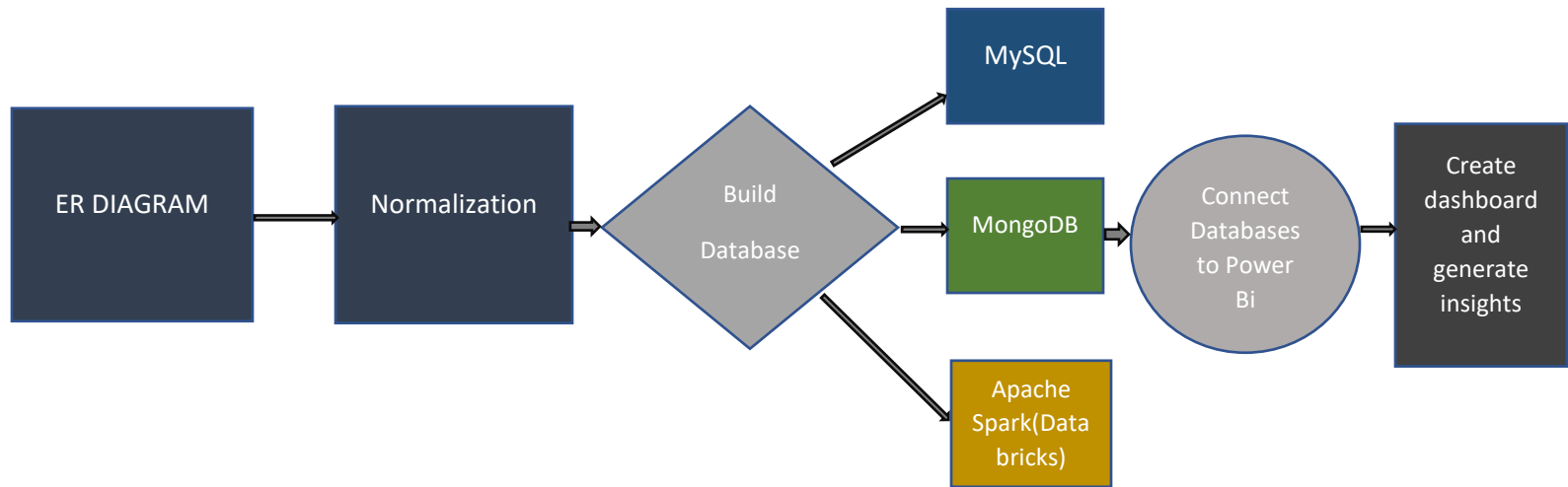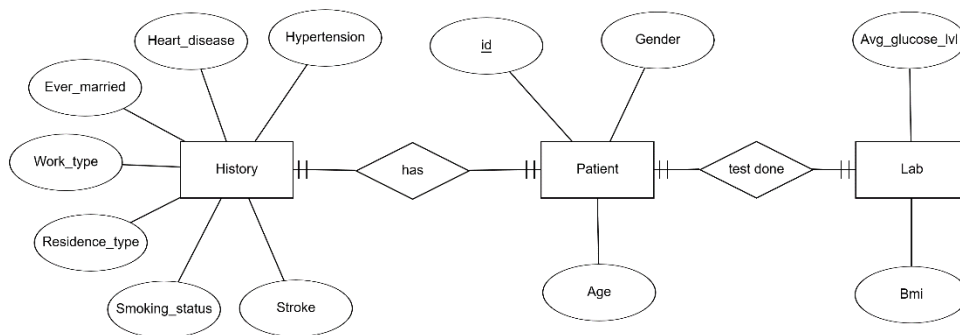# SYTEM ARCHITECTURE



## Conceptual Modeling with ER Diagram

The ER diagram would help us understand the entities, their attributes, and the relationships between them.

# Summary

In the ER diagram:

- Patient table contains personal attributes of an individual.
- Lifestyle table contains lifestyle-related attributes.
- Health table contains Health-related attributes.
- id is the primary key (PK) in the Patient table and a foreign key (FK) in the History table, and Lab table, linking these tables together.
- A Person entity has a one-to-one relationship with a history entity and a hospital entity, meaning that each person has one History and one Lab test results, and vice versa.

## Normalization

I normalized the data to **1NF**, **2NF**, and **3NF**.

**First Normal Form(1NF):**

- To achieve the first normal form, we need to make sure that each cell contains only a single value, each column has a unique name, and there are no duplicate rows. We also need to identify a primary key for the table. In this case, we can use the `id` column as the primary key since it uniquely identifies each row.
- No multivalued attributes are present.
- Based on the dataset, we already have a tabular representation, which satisfies 1NF.

**Second Normal Form (2NF)**:

- To achieve the second normal form, we need to make sure that each non-key attribute is fully dependent on the primary key. This means that there should be no partial dependencies, where a column depends on only part of the primary key. In this case, since the primary key is a single column (id), there are no partial dependencies
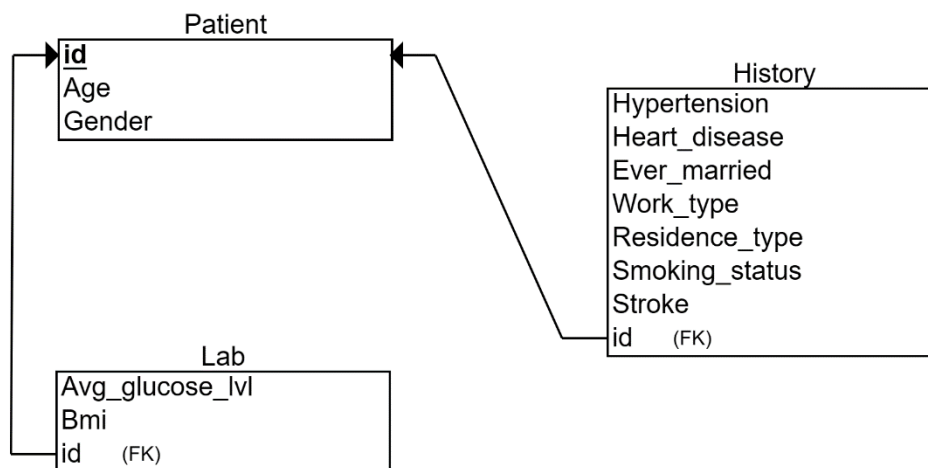- Therefore, the table already satisfies 2NF.

**Third Normal Form (3NF)**:
- To achieve the third normal form, we need to make sure that all non-key attributes are independent of each other. This means there should be no transitive dependencies, where a column depends on another column not the primary key.
- In summary, the existing table is already in 1NF and 2NF and 3NF

# Relational Schema

I translated the conceptual model into a relational schema, which includes defining tables, columns, primary keys, foreign keys, and other constraints. This step involves making decisions about how the data will be stored in a relational database.

The Relational Schema for the ER diagram is below:

### Patient
**id**
Age
Gender

### History
Hypertension
Heart_disease
Ever_married
Work_type
Residence_type
Smoking_status
Stroke
id      (FK)

### Lab
Avg_glucose_lvl
Bmi
id      (FK)

# Database Development using MySQL, MongoDB, and Apache Spark

## MySQL

### Objective: Develop a database using MySQL

**Database Implementation**

Using the Relational Schema. I partitioned the dataset into three distinct entities: Patient, History, and Lab. I employed MySQL to create three tables corresponding to these entities. Firstly, I created a new database 'health_stroke_data. I populated these tables by importing the relevant attributes into them from the dataset in .csv. Keeping them in separate tables would enhance the clarity and readability of the database design.
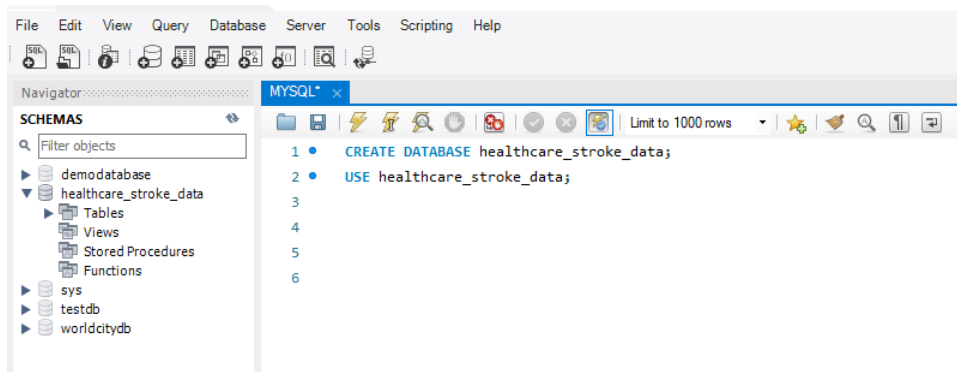
### Create Database
  ➢ Get Csv file

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_status | stroke | |
| 2 | 9046 | Male | 67 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 | formerly smoked | 1 | |
| 3 | 51676 | Female | 61 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | N/A | never smoked | 1 | |
| 4 | 31112 | Male | 80 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 | never smoked | 1 | |
| 5 | 60182 | Female | 49 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 | smokes | 1 | |
| 6 | 1665 | Female | 79 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24 | never smoked | 1 | |
| 7 | 56669 | Male | 81 | 0 | 0 | Yes | Private | Urban | 186.21 | 29 | formerly smoked | 1 | |
| 8 | 53882 | Male | 74 | 1 | 1 | Yes | Private | Rural | 70.09 | 27.4 | never smoked | 1 | |
| 9 | 10434 | Female | 69 | 0 | 0 | No | Private | Urban | 94.39 | 22.8 | never smoked | 1 | |
| 10 | 27419 | Female | 59 | 0 | 0 | Yes | Private | Rural | 76.15 | N/A | Unknown | 1 | |
| 11 | 60491 | Female | 78 | 0 | 0 | Yes | Private | Urban | 58.57 | 24.2 | Unknown | 1 | |
| 12 | 12109 | Female | 81 | 1 | 0 | Yes | Private | Rural | 80.43 | 29.7 | never smoked | 1 | |
| 13 | 12095 | Female | 61 | 0 | 1 | Yes | Govt_job | Rural | 120.46 | 36.8 | smokes | 1 | |
| 14 | 12175 | Female | 54 | 0 | 0 | Yes | Private | Urban | 104.51 | 27.3 | smokes | 1 | |
| 15 | 8213 | Male | 78 | 0 | 1 | Yes | Private | Urban | 219.84 | N/A | Unknown | 1 | |
| 16 | 5317 | Female | 79 | 0 | 1 | Yes | Private | Urban | 214.09 | 28.2 | never smoked | 1 | |
| 17 | 58202 | Female | 50 | 1 | 0 | Yes | Self-employed | Rural | 167.41 | 30.9 | never smoked | 1 | |
| 18 | 56112 | Male | 64 | 0 | 1 | Yes | Private | Urban | 191.61 | 37.5 | smokes | 1 | |
| 19 | 34120 | Male | 75 | 1 | 0 | Yes | Private | Urban | 221.29 | 25.8 | smokes | 1 | |
| 20 | 27458 | Female | 60 | 0 | 0 | No | Private | Urban | 89.22 | 37.8 | never smoked | 1 | |
| 21 | 25226 | Male | 57 | 0 | 1 | No | Govt_job | Urban | 217.08 | N/A | Unknown | 1 | |

➢ Create a new version of the .csv file without fieldnames (only data rows) By deleting the filed names

| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 9046 | Male | 67 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 | formerly smoked | 1 |
| 2 | 51676 | Female | 61 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | N/A | never smoked | 1 |
| 3 | 31112 | Male | 80 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 | never smoked | 1 |
| 4 | 60182 | Female | 49 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 | smokes | 1 |
| 5 | 1665 | Female | 79 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24 | never smoked | 1 |
| 6 | 56669 | Male | 81 | 0 | 0 | Yes | Private | Urban | 186.21 | 29 | formerly smoked | 1 |
| 7 | 53882 | Male | 74 | 1 | 1 | Yes | Private | Rural | 70.09 | 27.4 | never smoked | 1 |
| 8 | 10434 | Female | 69 | 0 | 0 | No | Private | Urban | 94.39 | 22.8 | never smoked | 1 |
| 9 | 27419 | Female | 59 | 0 | 0 | Yes | Private | Rural | 76.15 | N/A | Unknown | 1 |
| 10 | 60491 | Female | 78 | 0 | 0 | Yes | Private | Urban | 58.57 | 24.2 | Unknown | 1 |
| 11 | 12109 | Female | 81 | 1 | 0 | Yes | Private | Rural | 80.43 | 29.7 | never smoked | 1 |
| 12 | 12095 | Female | 61 | 0 | 1 | Yes | Govt_job | Rural | 120.46 | 36.8 | smokes | 1 |
| 13 | 12175 | Female | 54 | 0 | 0 | Yes | Private | Urban | 104.51 | 27.3 | smokes | 1 |
| 14 | 8213 | Male | 78 | 0 | 1 | Yes | Private | Urban | 219.84 | N/A | Unknown | 1 |
| 15 | 5317 | Female | 79 | 0 | 1 | Yes | Private | Urban | 214.09 | 28.2 | never smoked | 1 |
| 16 | 58202 | Female | 50 | 1 | 0 | Yes | Self-employed | Rural | 167.41 | 30.9 | never smoked | 1 |
| 17 | 56112 | Male | 64 | 0 | 1 | Yes | Private | Urban | 191.61 | 37.5 | smokes | 1 |
| 18 | 34120 | Male | 75 | 1 | 0 | Yes | Private | Urban | 221.29 | 25.8 | smokes | 1 |
| 19 | 27458 | Female | 60 | 0 | 0 | No | Private | Urban | 89.22 | 37.8 | never smoked | 1 |
| 20 | 25226 | Male | 57 | 0 | 1 | No | Govt_job | Urban | 217.08 | N/A | Unknown | 1 |
| 21 | 70630 | Female | 71 | 0 | 0 | Yes | Govt_job | Rural | 193.94 | 22.4 | smokes | 1 |

➢ Create database and switch to new database

```
CREATE DATABASE healthcare_stroke_data;
USE healthcare_stroke_data;
```
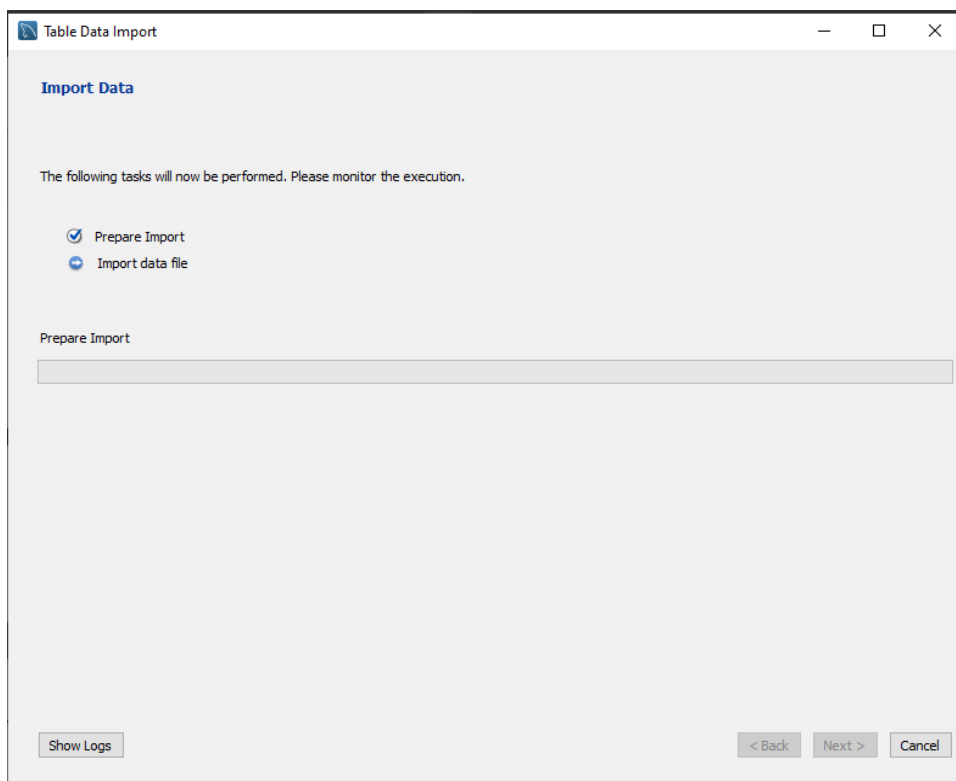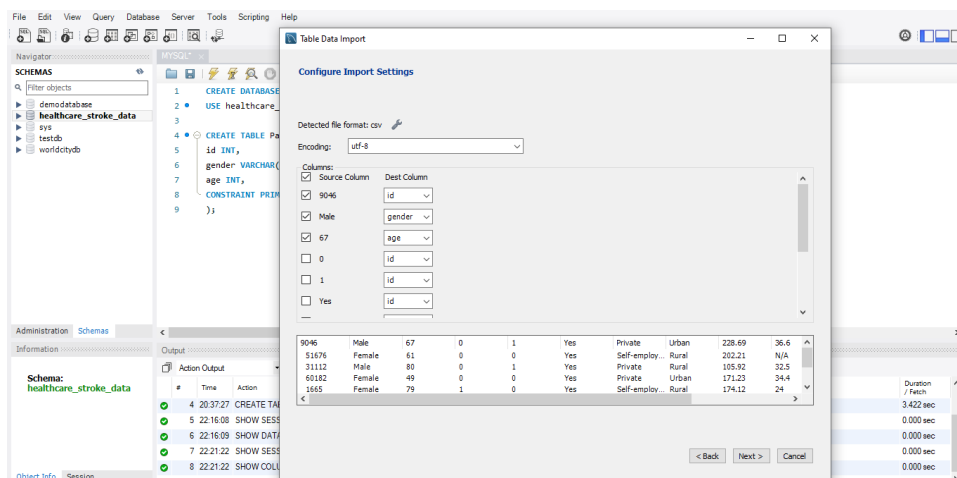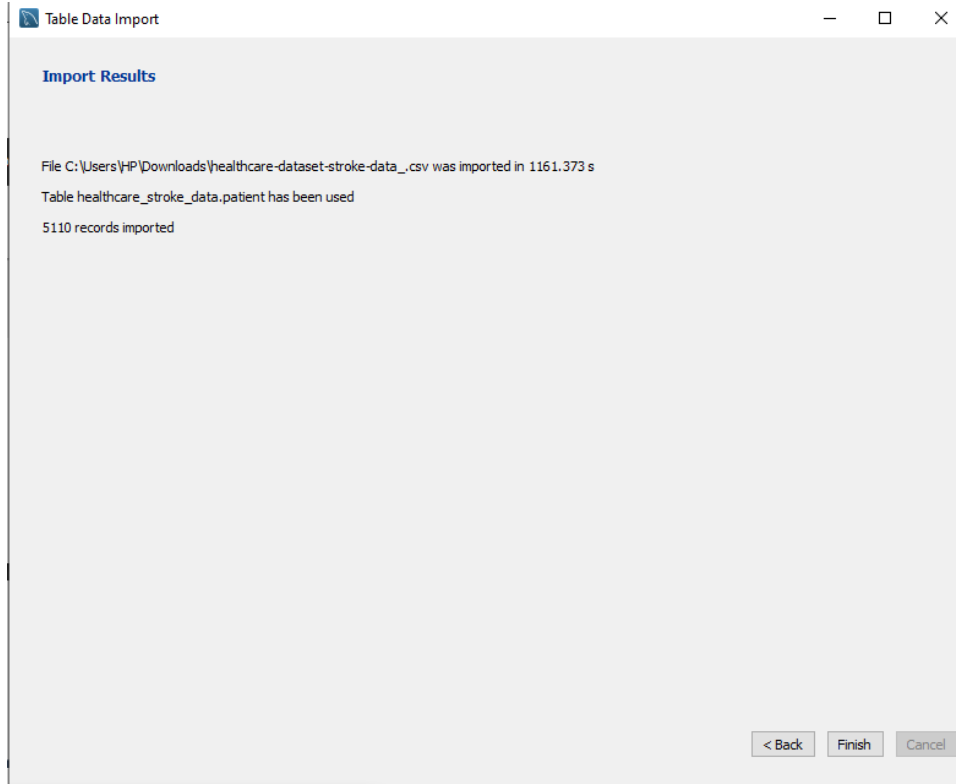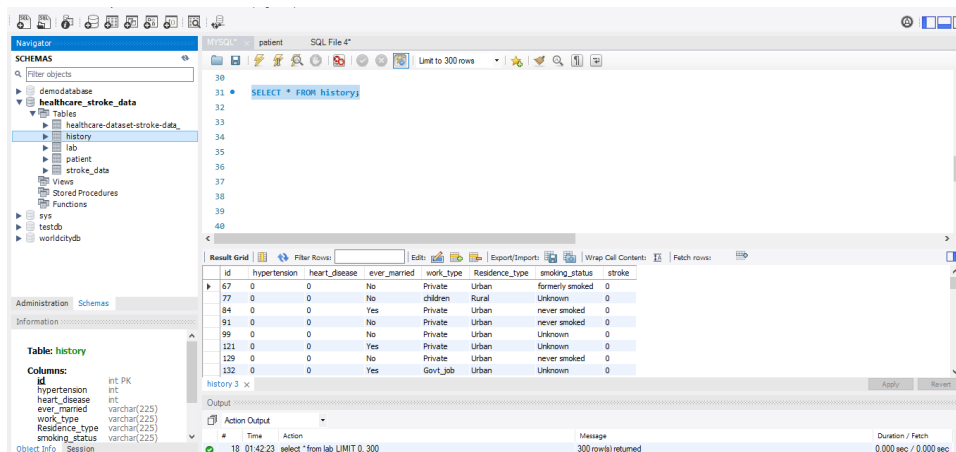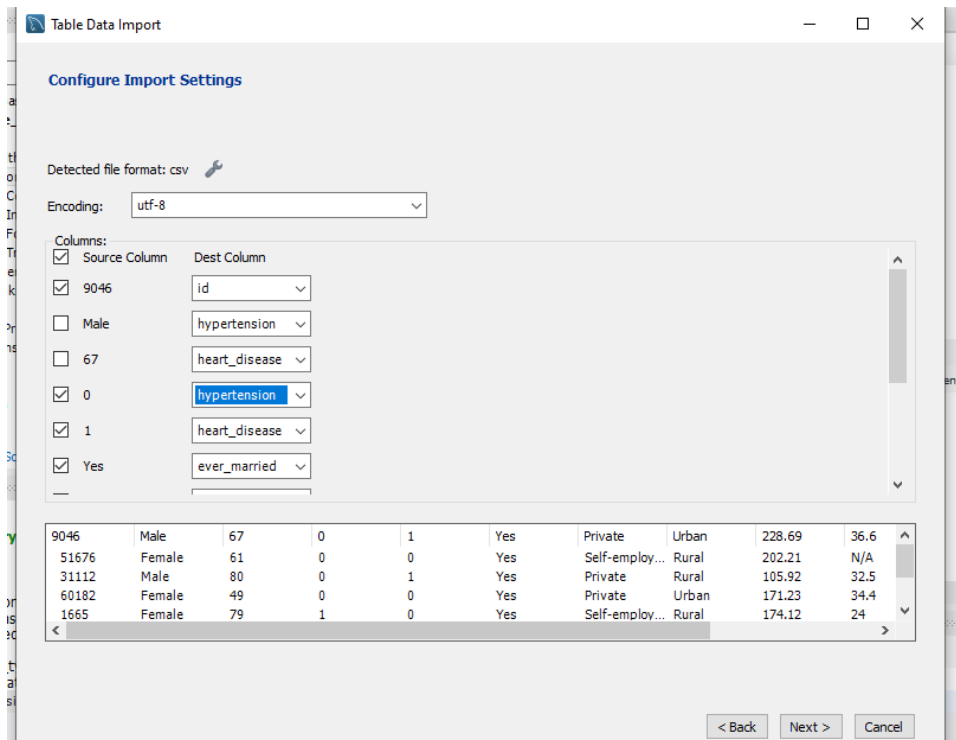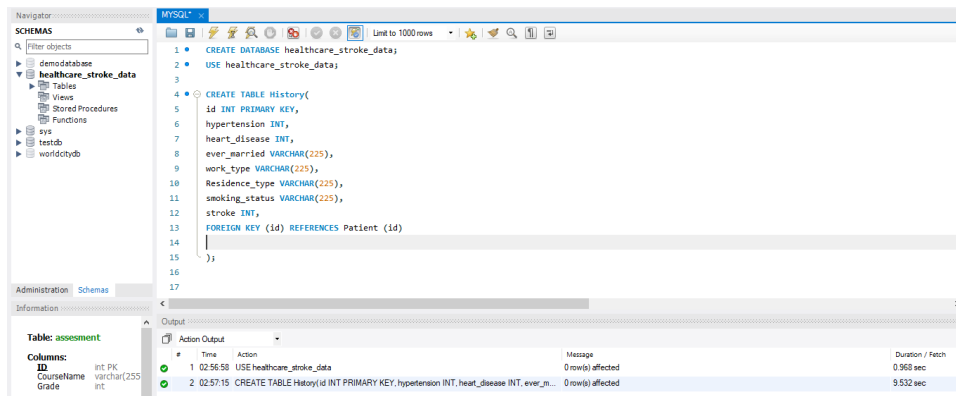
➢ Create table "Patient"

> ➤ Imported Csv file to the table

Table Data Import — Select Destination

Select destination table and additional options.

- Use existing table: healthcare_stroke_data.patient
- Create new table: healthcare_stroke_dat ▾ healthcare-dataset-stroke-data_
- Drop table if exists



Table Data Import — Configure Import Settings

Detected file format: csv

Encoding: utf-8

Columns:

| Source Column | Dest Column |
|---|---|
| 9046 | id |
| Male | gender |
| 67 | age |
| 0 | id |
| 1 | id |
| Yes | id |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 9046 | Male | 67 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 |
| 51676 | Female | 61 | 0 | 0 | Yes | Self-employ... | Rural | 202.21 | N/A |
| 31112 | Male | 80 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 |
| 60182 | Female | 49 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 |
| 1665 | Female | 79 | 1 | 0 | Yes | Self-employ... | Rural | 174.12 | 24 |

< Back    Next >    Cancel



Table Data Import

**Import Data**

The following tasks will now be performed. Please monitor the execution.

- ✔ Prepare Import
- ● Import data file

Prepare Import

Show Logs                     < Back    Next >    Cancel

**Table Data Import**

**Import Results**

File C:\Users\HP\Downloads\healthcare-dataset-stroke-data_.csv was imported in 1161.373 s

Table healthcare_stroke_data.patient has been used

5110 records imported

‹ Back | Finish | Cancel

➢ Queried the data



```
1   CREATE DATABASE healthcare_stroke_data;
2   USE healthcare_stroke_data;
3
4   CREATE TABLE Patient(
5       id INT,
6       gender VARCHAR(7),
7       age INT,
8       CONSTRAINT PRIMARY KEY (id)
9   );
10
11  SELECT * FROM Patient
12  ORDER BY age DESC;
```

I repeated the process for other tables like "History" and "Lab" Entities

➢ **For History**

➢ **For Lab**

The three entities have been successfully portioned into three tables: Patient, History and Lab.

## Testing and Optimization

Query optimization is a crucial aspect of database management aimed at improving the performance and efficiency of SQL queries.



In this query:

- Not all columns are selected: Instead of selecting all columns using SELECT *, explicitly specify only the required columns in the SELECT statement. This reduces the amount of data transferred and can improve query performance.
- JOIN conditions are efficient and utilize indexes
- WHERE clauses filter data effectively
- The result set is limited to 10 rows: Use the LIMIT clause to restrict the number of rows returned by a query, especially for queries that return large result sets.

By following these optimization techniques and continuously fine-tuning queries based on performance analysis, we can significantly improve the efficiency and responsiveness of our database.

## NoSQL

### Objective: Develop a database using MongoDB
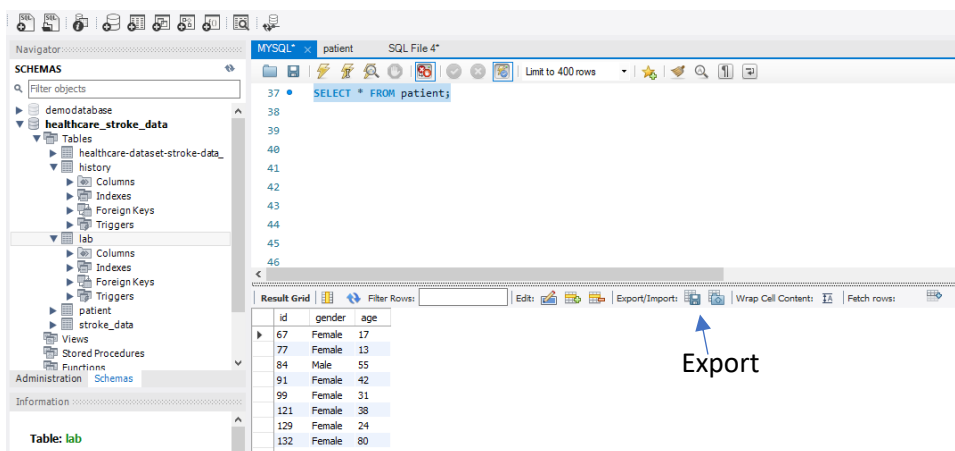To develop a MongoDB database. The following steps will be taken.

➢ **Collections**

From the ER diagram, the dataset structure and relationships was known. With this I could define the collections which are: Patients, History and Lab.
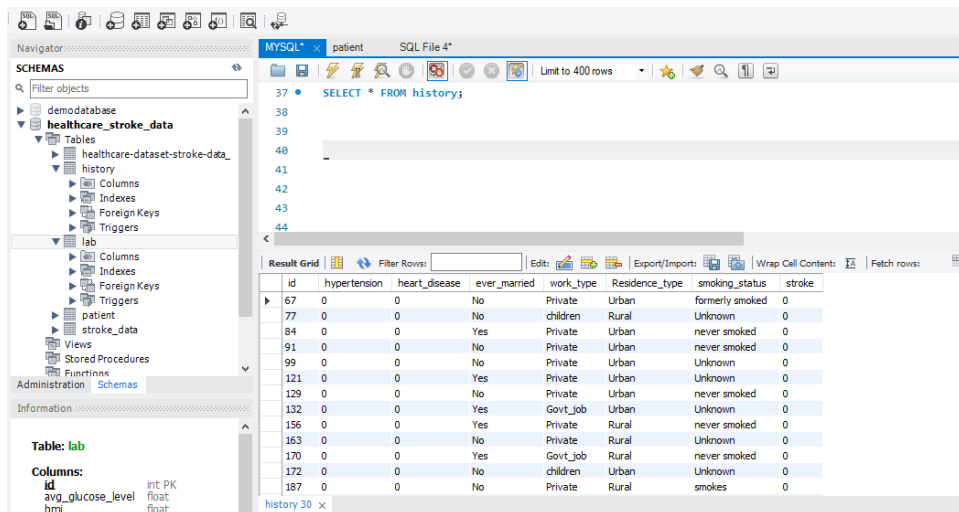
### Create Database
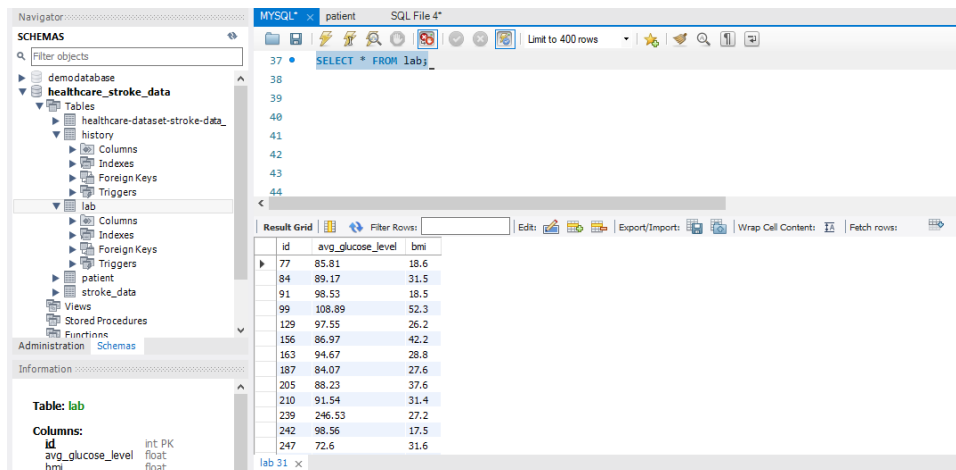I migrated the dataset from MySQL to MongoDB

➢ Imported data from MySQL Workbench



➢ Did the same for the other two entities

- ➢ Opened MondoDB Compass and created a database named healthcare_stroke_data with three collections (Patient, History and Lab)



 history.csv

 patient.csv

 lab.csv

The data was imported as .csv successfully.

- ➢ Created a new database "health_stroke_data "and three collections: Patient, History and Lab.

➢ Inserted Documents or Add Data

Added data to the collections by importing csv file from fig

Data was successfully imported to the patient collection. I followed the same steps for the other two collections.

➢ Queried the Data

I used MongoDB's query language to retrieve data from the database.

```
>_MONGOSH
> use health_stroke_data
< switched to db health_stroke_data
> // Find all patients
  db.patient.find()

< {
    _id: ObjectId('660603d81a10b916acad233c'),
    id: 67,
    gender: 'Female',
    age: 17
  }
  {
    _id: ObjectId('660603d81a10b916acad233d'),
    id: 77,
    gender: 'Female',
    age: 13
  }
  {
    _id: ObjectId('660603d81a10b916acad233e'),
    id: 84,
    gender: 'Male',
    age: 55
```
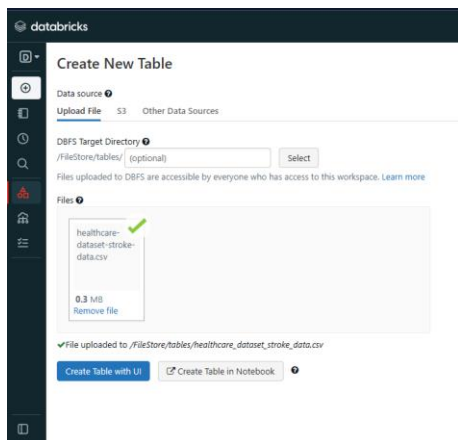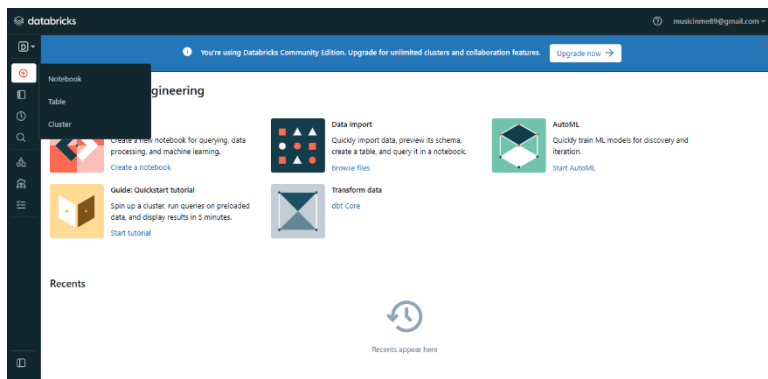
```
>_MONGOSH
    smoking_staus: 'smokes'
  }
  Type "it" for more
> // To find patients that are "married" and "smokes" and diagnosed with "stroke"
  db.history.find({}, {ever_married:'Yes', smoking_staus:'smokes', stroke:'1'})
< {
    _id: ObjectId('660604121a10b916acad24d5'),
    ever_married: 'Yes',
    smoking_staus: 'smokes',
    stroke: '1'
  }
  {
    _id: ObjectId('660604121a10b916acad24cf'),
    ever_married: 'Yes',
    smoking_staus: 'smokes',
    stroke: '1'
  }
  {
    _id: ObjectId('660604121a10b916acad24dd'),
    ever_married: 'Yes',
    smoking_staus: 'smokes',
    stroke: '1'
  }
```

## Apache Spark

### Objective: Develop a database using Data Bricks

To develop a MongoDB database. The following steps will be taken.

➢ Uploaded the csv file

The csv file was uploaded successfully.

▾ (3) Spark Jobs

    ▸ Job 4   View (Stages: 1/1)
    ▸ Job 5   View (Stages: 1/1)
    ▸ Job 6   View (Stages: 1/1)

▾ ▦ df: pyspark.sql.dataframe.DataFrame

```
id: integer
gender: string
age: double
hypertension: integer
heart_disease: integer
ever_married: string
work_type: string
Residence_type: string
avg_glucose_level: double
bmi: string
smoking_status: string
stroke: integer
```

The image above shows the data frame and the datatypes of each column the dataset uploaded.

➢ Queried the data
➢ Created a temp table

Cmd 3

```python
1   # Create a view or table
2
3   temp_table_name = "healthcare_dataset_stroke_data_csv2"
4
5   df.createOrReplaceTempView(temp_table_name)
```

Command took 0.12 seconds -- by musicinme89@gmail.com at 4/2/2024, 4:48:05 PM on Healthcare-dataset-stroke-data

*Figure 1.1 1*

Cmd 4

```sql
1   %sql
2
3   /* Query the created temp table in a SQL cell */
4
5   select * from `healthcare_dataset_stroke_data_csv2`
```

▸ (1) Spark Jobs
▸ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [id: integer, gender: string ... 10 more fields]

Table ∨ +

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type | avg_glucose_level | bmi | smoking_statu |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 9046 | Male | 67 | 0 | 1 | Yes | Private | Urban | 228.69 | 36.6 | formerly smok |
| 2 | 51676 | Female | 61 | 0 | 0 | Yes | Self-employed | Rural | 202.21 | N/A | never smoked |
| 3 | 31112 | Male | 80 | 0 | 1 | Yes | Private | Rural | 105.92 | 32.5 | never smoked |
| 4 | 60182 | Female | 49 | 0 | 0 | Yes | Private | Urban | 171.23 | 34.4 | smokes |
| 5 | 1665 | Female | 79 | 1 | 0 | Yes | Self-employed | Rural | 174.12 | 24 | never smoked |
| 6 | 56669 | Male | 81 | 0 | 0 | Yes | Private | Urban | 186.21 | 29 | formerly smok |
| 7 | 53882 | Male | 74 | 1 | 1 | Yes | Private | Rural | 70.09 | 27.4 | never smoked |

⬇ 5,110 rows | 0.76 seconds runtime    Refreshed now

ⓘ SQL cell result stored as PySpark data frame _sqldf . Learn more

Command took 0.76 seconds -- by musicinme89@gmail.com at 4/2/2024, 4:49:08 PM on Healthcare-dataset-stroke-data

➢ Visualized trends and insights.

From the fig above. People who previously smoked were at a higher risk of being diagnosed with stroke.

➢ Saved the temp table



```python
# With this registered as a temp view, it will only be available to this particular notebook. If you'd like other users to be able to query this table, you can also
create a table from the DataFrame.
# Once saved, this table will persist across cluster restarts as well as allow various users across different notebooks to query this data.
# To do so, choose your table name and uncomment the bottom line.

permanent_table_name = "healthcare_dataset_stroke_data_csv2"

# df.write.format("parquet").saveAsTable(permanent_table_name)
```

Command took 0.10 seconds -- by musicinme89@gmail.com at 4/2/2024, 5:15:21 PM on Healthcare-dataset-stroke-data

# DATA VISUALIZATION ON POWER BI

I integrated MySQL with Power BI to create an interactive dashboard and report to generate meaningful insights from the data.

➢ Connected Power BI to SQL Server database



➢ Got the server from MySQL

**MySQL Connections** ⊕ ⊗

Local instance MySQL80
👤 root
localhost:3306

**MySQL database**                                              ✕

Server
127.0.0.1:3306

Database
healthcare_stroke_data

▷ Advanced options

                                              OK        Cancel

---

MySql database                                              ✕

🗄 127.0.0.1:3306;healthcare_stroke_data

**Windows**        Use your Windows credentials to access this database.

**Database**       ⦿ Use my current credentials
                   ○ Use alternate credentials
                   User name

                   Password

                   Select which level to apply these settings to
                   127.0.0.1:3306                              ▾

        Back                              Connect    Cancel

---

➢ Selected the relevant tables from the database (Patient, History, Lab)

**Navigator**                                              ☐ ✕

[Search]  🔍                    healthcare_stroke_data.patient

Display Options ▾              | id | gender | age | healthcare_stroke_data.history | healthcare_stroke_ |
◢ 🗄 127.0.0.1:3306: healthcare_stroke_data [5]  | 67 | Female | 17 | Table | Table |
  ☐ ▦ healthcare_stroke_data.healthcare-dataset-...  | 77 | Female | 13 | Table | Table |
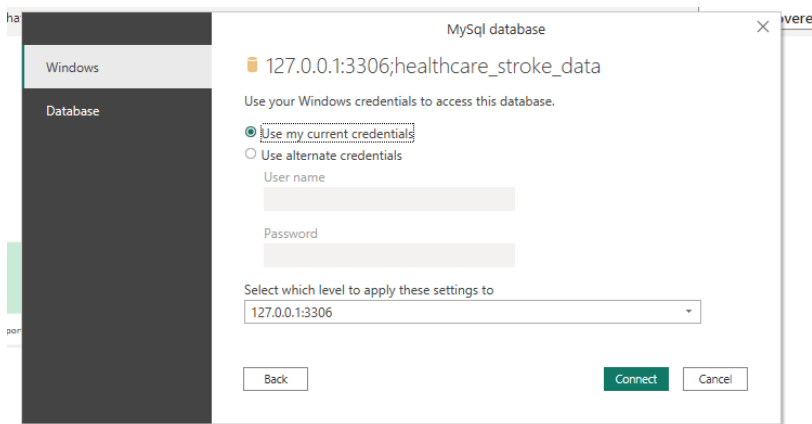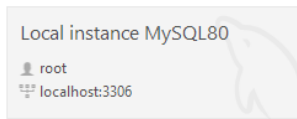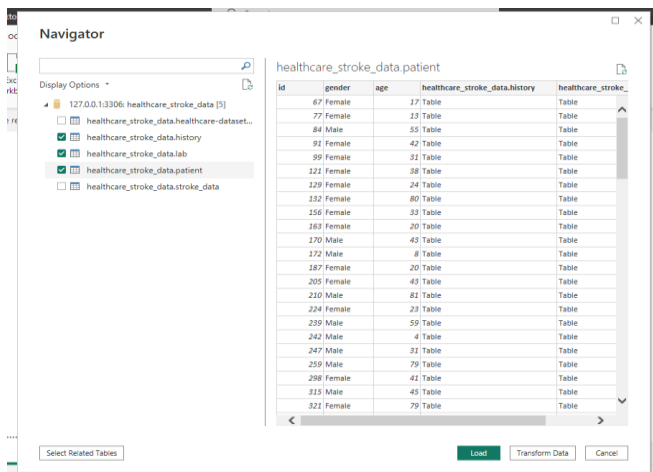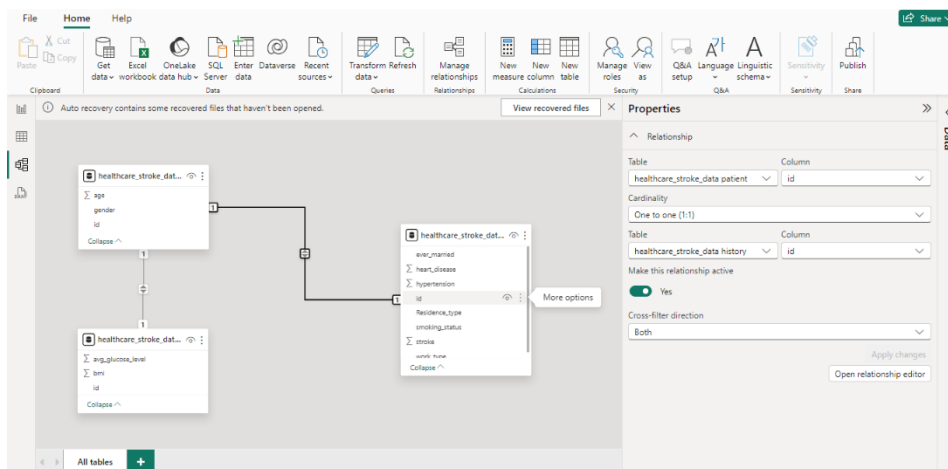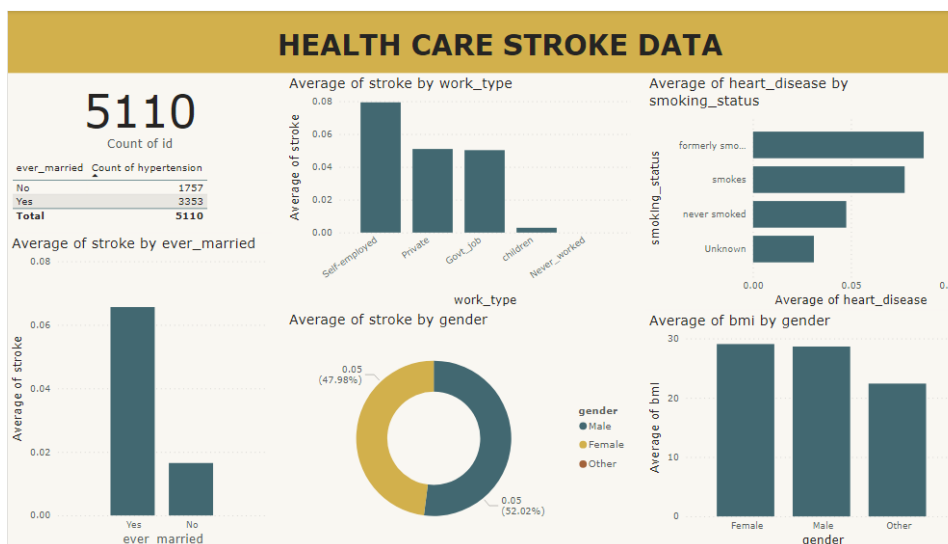  ☑ ▦ healthcare_stroke_data.history  | 84 | Male | 55 | Table | Table |
  ☑ ▦ healthcare_stroke_data.lab  | 91 | Female | 42 | Table | Table |
  ☑ ▦ healthcare_stroke_data.patient  | 99 | Female | 31 | Table | Table |
  ☐ ▦ healthcare_stroke_data.stroke_data  | 121 | Female | 38 | Table | Table |
                               | 129 | Female | 24 | Table | Table |
                               | 132 | Female | 80 | Table | Table |
                               | 156 | Female | 33 | Table | Table |
                               | 163 | Female | 20 | Table | Table |
                               | 170 | Male | 43 | Table | Table |
                               | 172 | Male | 8 | Table | Table |
                               | 187 | Female | 20 | Table | Table |
                               | 205 | Female | 43 | Table | Table |
                               | 210 | Male | 81 | Table | Table |
                               | 224 | Female | 23 | Table | Table |
                               | 239 | Male | 59 | Table | Table |
                               | 242 | Male | 4 | Table | Table |
                               | 247 | Male | 31 | Table | Table |
                               | 259 | Male | 79 | Table | Table |
                               | 298 | Female | 41 | Table | Table |
                               | 315 | Male | 45 | Table | Table |
                               | 321 | Female | 79 | Table | Table |

Select Related Tables                    Load    Transform Data    Cancel

---

➢ Created Relationships

➢ Created an interactive dashboards and report



# Insights derived from the dashboard data:

- We have 5110 patients in our dataset.
- Patients with a history of marriage tend to have a higher incidence of hypertension.
- Self-employed individuals show a higher propensity for stroke compared to other occupational groups, while children were not diagnosed with stroke.
- Heart diseases are more prevalent among patients with a history of smoking.
- Patients with a marital history are more prone to being diagnosed with stroke.
- Men constitute 52.02% of stroke diagnoses, indicating a higher likelihood of stroke diagnosis in men compared to women that constitutes 47.98%.
- On average, women have a higher BMI than men.

# RECOMMENDATION

Based on the insights provided by the dashboard data, several recommendations can be made:

1. Targeted Screening and Education: Focus on screening and education for hypertension among married individuals.

2. Stroke Prevention for Self-Employed: Develop strategies to prevent stroke among self-employed individuals, including stress management and awareness campaigns.

3. Smoking Cessation Programs: Implement smoking cessation initiatives to reduce heart disease among smokers.

4. Stroke Awareness Campaigns: Launch campaigns to raise awareness about stroke prevention and symptoms, particularly targeting married individuals.

5. Gender-Specific Health Initiatives: Develop health programs tailored to address the higher stroke risk in men and higher BMI in women.

6. Occupational Health Interventions: Provide occupational health support for self-employed individuals to mitigate stroke risk factors related to their work.