# Predicting stroke probability

T5 BOOTCAMP DATA SCIENCE PROJECT

# Outlines

- Project goal

- Dataset

- Data cleaning

- Data processing

- Visualize data

- Models

# Project Goal

The goal of this project is to use classification model to predict which patients are close to have stroke and compute the accuracy of this model. These outcomes will help to focus on potential patients and avoid them being in stroke critical case.

# Dataset

- Data provided by Kaggle.com has been used in this project.

- The dataset contains 5110 unique records with 12 attributes for each, collecting from 2995 females and 2115 males. The last column contains '1' if the patient had stroke and '0' if he or she hadn't. A few attributes like age, marriage, and hypertension will be consider in prediction model. The dataset divided in three parts:

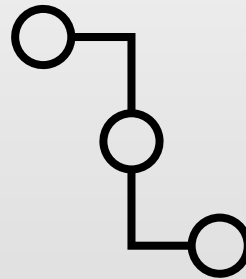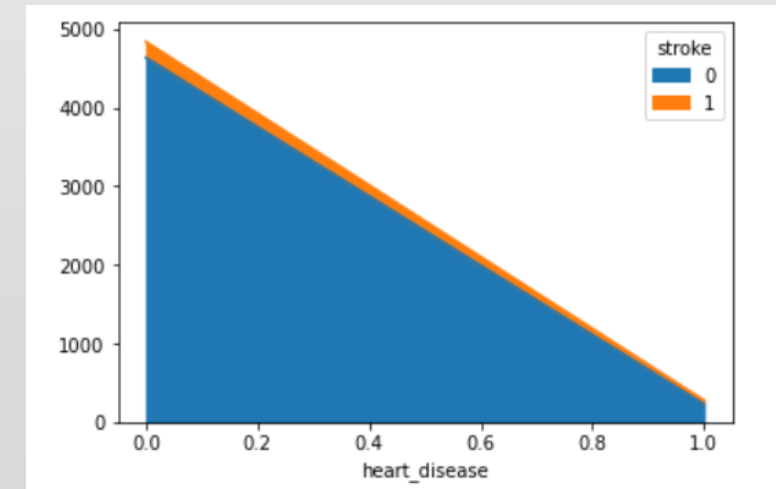- 20% for testing.

- 60% for training.
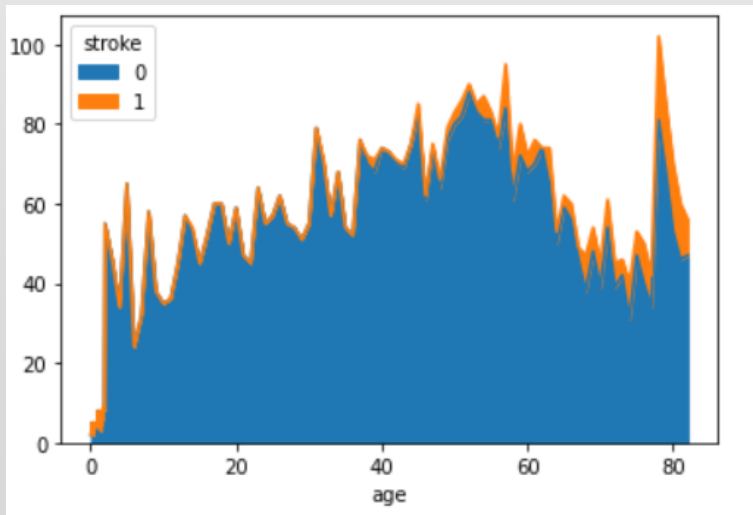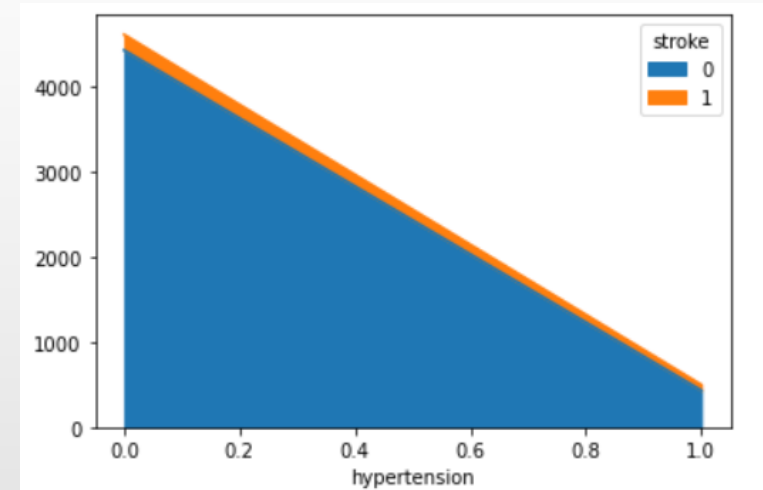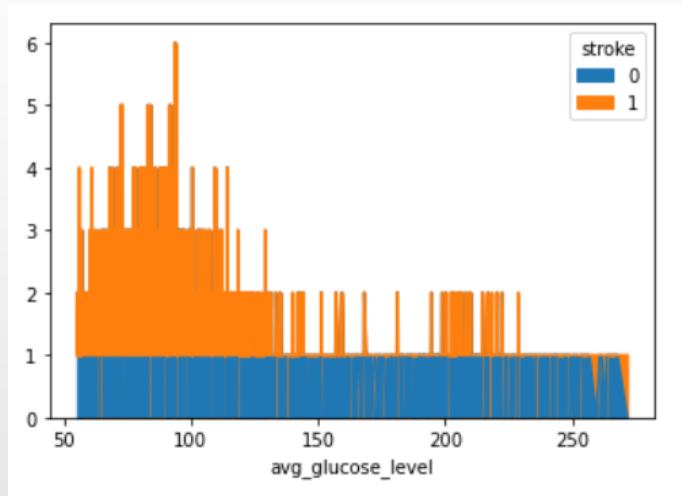
- 20% for validation.

# Data Cleaning

- Missing Data
  - Replaced with mean value

- Irregular Data (Outliers)
  - Replaced with mode value

- Unnecessary Data
  - Deleted it
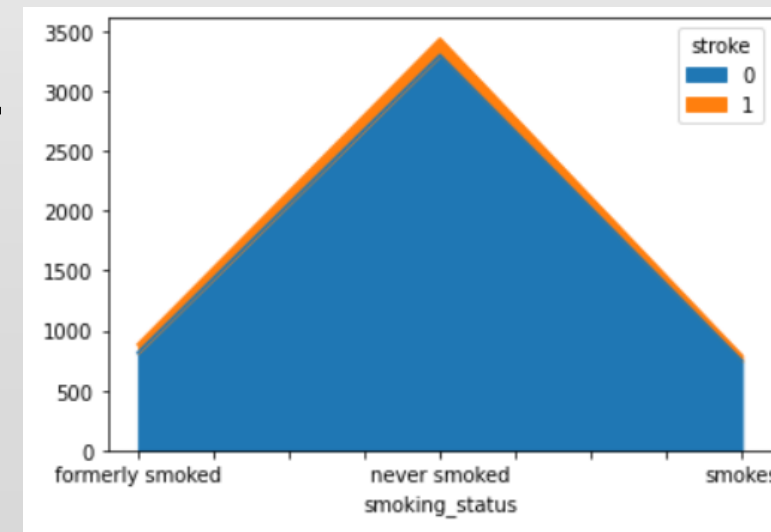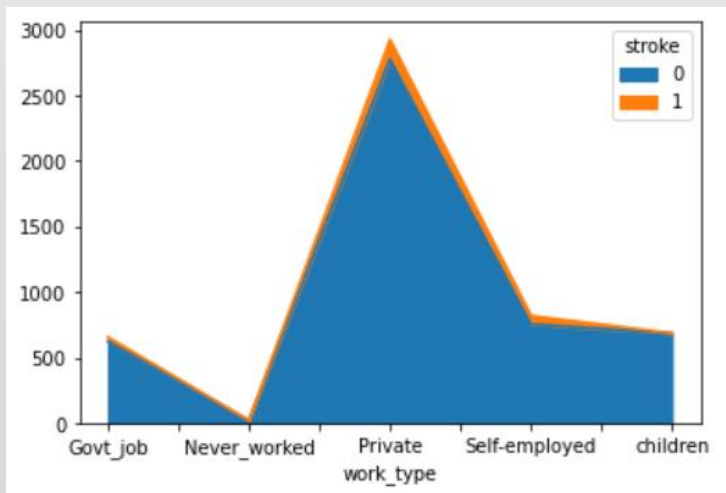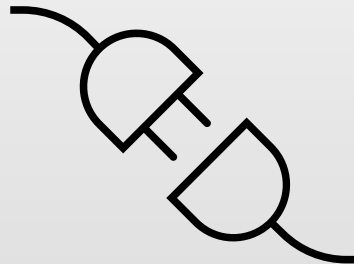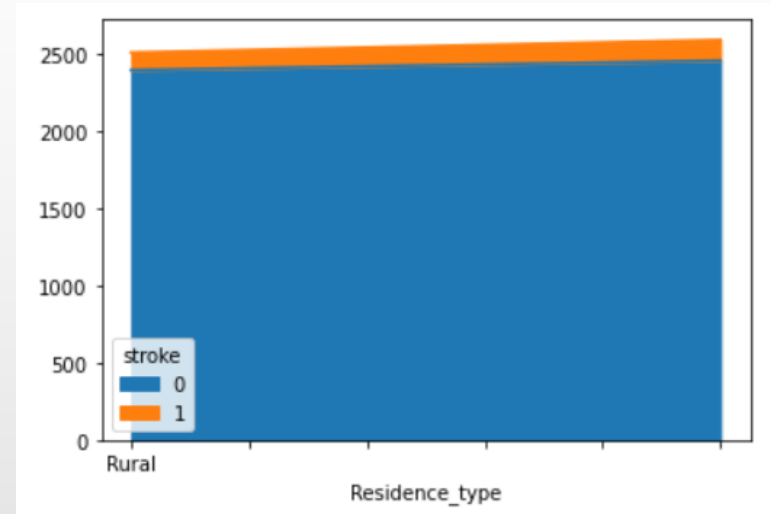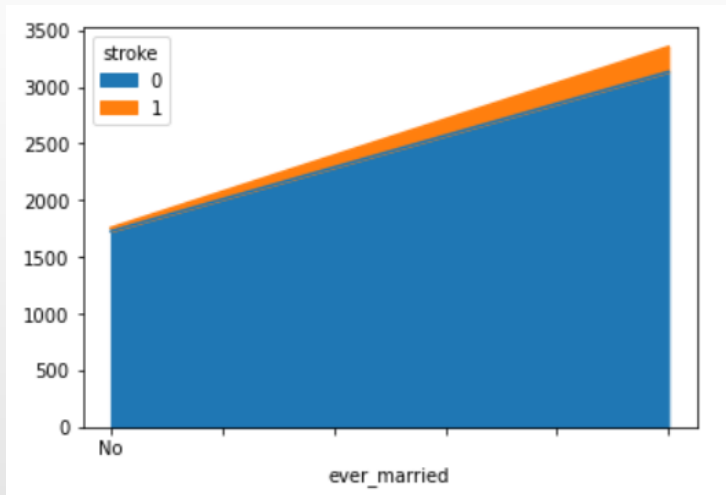
- Inconsistent Data
  - Nothing

# Data Processing

- Convert categorical variables to numbers by LabelEncoder from sklearn.

- Scale values by using StandardScaler from sklearn.
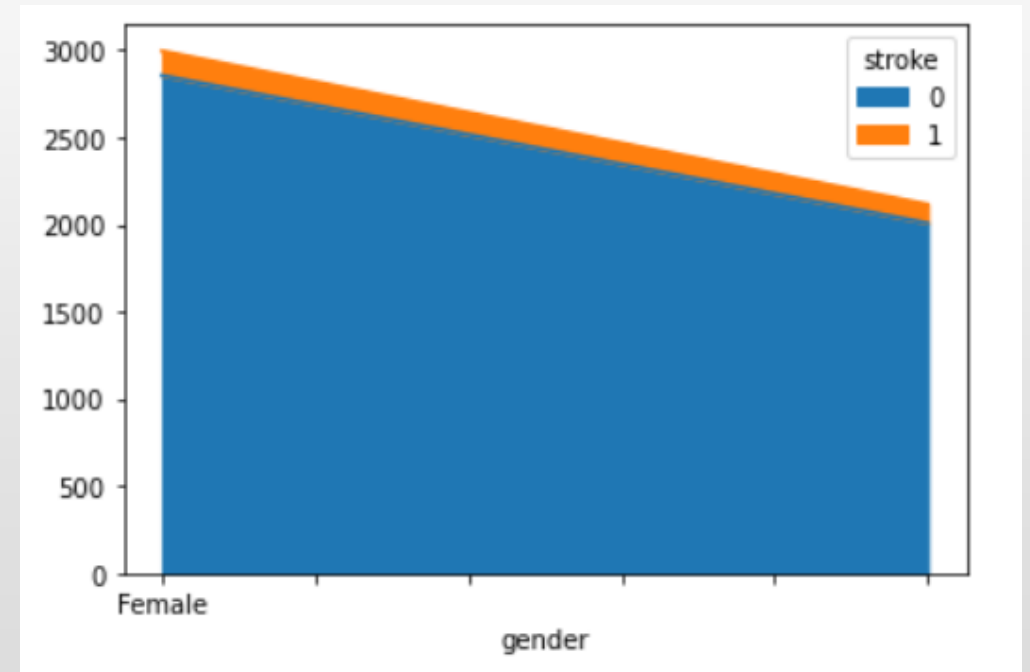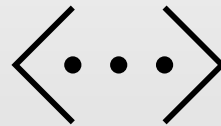
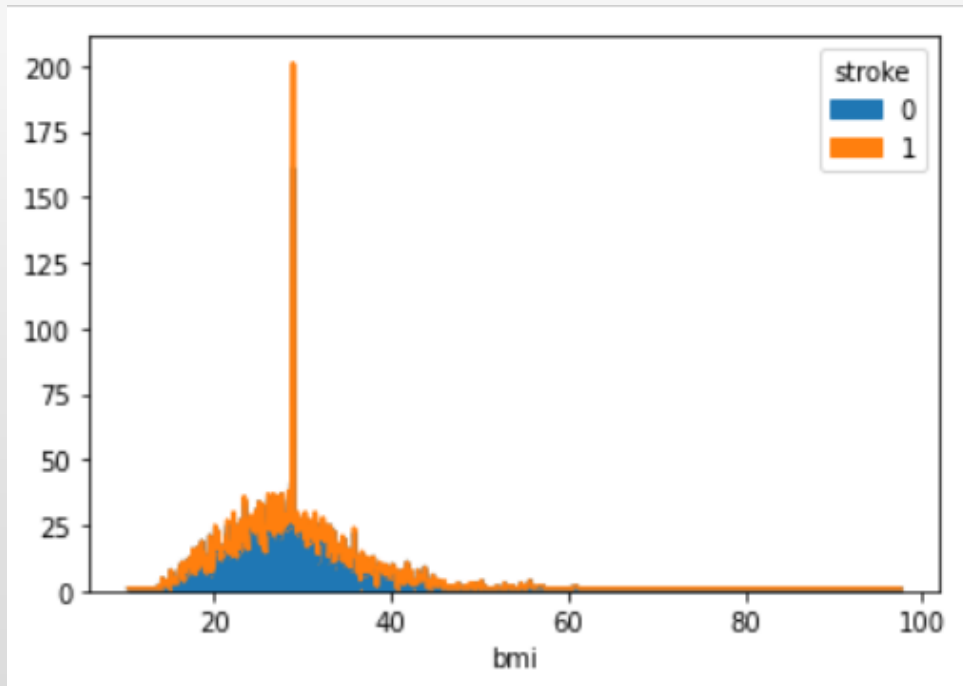- Use SMOTE from imblearn to solve the imbalence data
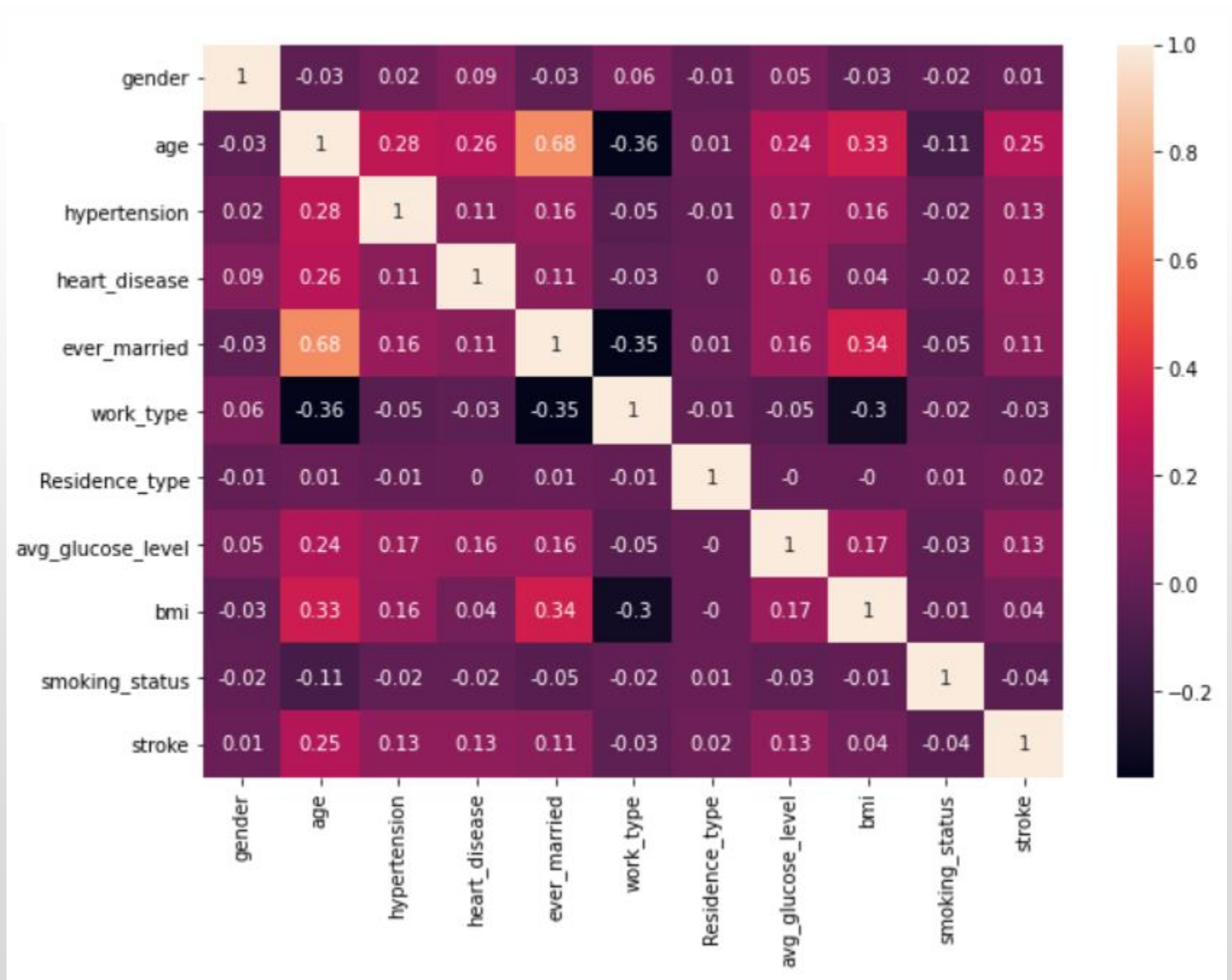
# Data Visualization

# Data Visualization

# Data Visualization

All together

# Models

- Models' accuracy

- Logistic Regression: 0.786692759295499

- Random Forest: 0.910958904109589

- Decision Tree: 0.8688845401174168

- KNN: 0.8140900195694716

- Naive Bayes: 0.8140900195694716

- KMeans: 0.8551859099804305

# Thank You