

Theoretical Machine Learning Project

Heart Stroke Prediction

By -

Thejaswin S(RA1911026010029)

Sarthak Gupta(RA1911026010041)

Shravya Sharan(RA1911026010055)

K1 SECTION

1. ABSTRACT

In recent times, Heart Stroke prediction has been one of the most complicated tasks in the medical field. In the modern era, approximately one person dies per minute due to a heart Stroke. Data science plays a crucial role in processing a massive amount of data in the field of healthcare. As heart stroke prediction is a complex task, there is a need to automate the prediction process to avoid associated risks and alert the patient well in advance.

Heart disease and strokes have rapidly increased globally, even at young ages. Stroke prediction is a complex task requiring a tremendous amount of data pre-processing. There is a need to automate the prediction process for the early detection of symptoms related to stroke so that it can be prevented at an early stage. In the proposed model, heart stroke prediction is performed on a dataset. We used machine learning algorithms like Gradient Boosting Classifier and Random Forest Classifier to train two different models for accurate prediction. The algorithm that best performed this task is Random Forest Classifier that gave an accuracy of approximately 90%.

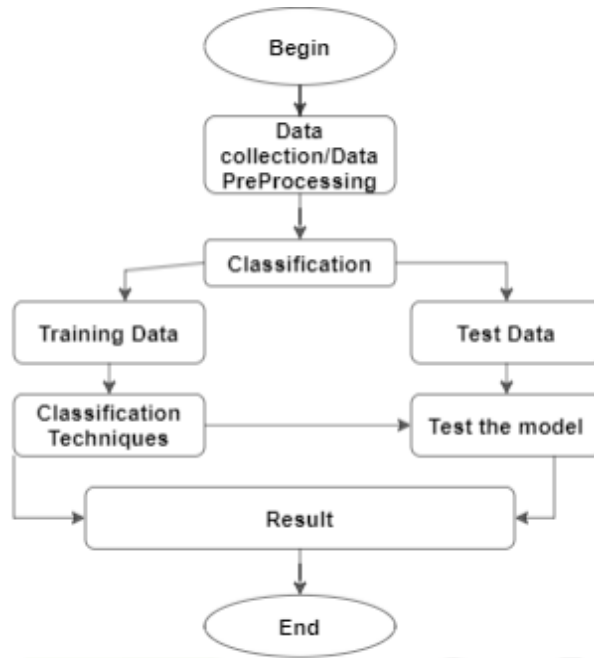
2. INTRODUCTION

Cardiovascular Diseases (CVDs) are the most common cause of death globally, representing 32% of all global deaths, with about 17.9 million people being affected. Out of these, the two most common CVDs are heart attack and heart strokes, accounting for 85% of the total people. Heart attack is caused due to blockage of oxygen or blood supply to the heart muscle, while heart stroke is caused when there is blockage of the vessel feeding the brain. Although both diseases are different from each other, the risk factors contributing to them are pretty similar. The risk factors include unhealthy diet, tobacco use, diabetes, sedentary lifestyle, harmful alcohol use, high blood pressure and family history. Detecting heart stroke and taking medical action immediately can prolong life and help prevent heart disease in the future.

Machine learning has become one of the most demanding fields in modern technology. It is a form of artificial intelligence where the model can analyse the data, identify patterns and predict the outcome with minimal human intervention. Various machine learning algorithms can make heart stroke predictions in adults. It has become an intriguing research problem as multiple factors or parameters can influence the outcome. The factors include work type, gender, residence type, age, average glucose level, body mass index, smoking status of the individual and any previous heart disease.

3. PROPOSED SYSTEM

The proposed work predicts heart stroke by exploring the above mentioned four classification algorithms and does performance analysis. The objective of this study is to predict if the patient suffers from heart stroke effectively. The health professional enters the input values from the patient's health report. The data is fed into the model, which predicts the probability of having a heart stroke. Fig. shows the entire process involved.



4. METHODOLOGY

This section is divided into two parts; these are Data description and Machine learning classifiers. These two processes are described below:

- a. *Data Description:* We used the heart stroke dataset available on the Kaggle website for our analysis. This dataset consists of a total of 12 attributes. The complete description of the attributes used in the proposed work is given below:
 - id: This attribute means a person's id. It's numerical data.
 - Age: This attribute means a person's age. It's numerical data.
 - Gender: This attribute means a person's gender. It's categorical data.
 - Hypertension: This attribute means that this person is hypertensive or not. It's numerical data.
 - Work type: This attribute represents the person work scenario. It's categorical data.
 - Residence type: This attribute represents the person living scenario. It's categorical data.
 - Heart disease: This attribute means whether this person has a heart disease person or not. It's numerical data.
 - Avg glucose level: This attribute means what was the level of a person's glucose condition. It's numerical data.
 - Bmi: This attribute means the body mass index of a person. It's numerical data.
 - Ever married: This attribute represents a person's married status. It's categorical data.
 - Smoking Status: This attribute means a person's smoking condition. It's categorical data.
 - Stroke: This attribute means a person previously had a stroke or not. It's numerical data.
 Attribute stroke is the decision class, and the rest of the attribute is the response class.
- b. *Machine Learning Classifiers:*
 - Random Forest classifier:
Random Forest algorithms are used for classification as well as regression. It creates a tree for the data and makes a prediction based on that. Random Forest algorithm can be used on large datasets and produce the same result even when large sets of record values are missing. The generated samples from the decision tree can be saved to be used on other data. There are two

stages in a random forest: firstly, create a random forest, then make a prediction using a random forest classifier built in the first stage.

- Gradient Boosting classifier:
Gradient boosting algorithm can be used for predicting continuous target variable (as a Regressor) and categorical target variable (as a Classifier). When used as a regressor, the cost function is Mean Square Error (MSE) and when it is used as a classifier, the cost function is Log loss.

5. IMPLEMENTATION

(code attached at the end of the report)

×

Select your gender

☒ Male

☐ Female

Do you have hypertension?

☒ Yes

☐ No

Do you have heart disease?

☒ Yes

☐ No

Have you ever been married?

☒ Yes

☐ No

Heart ♥ Stroke Prediction App

Select your residence type

Urban

How old are you?

879168

Select the type of work you do.

Govt job

How much is your bmi?

1882168

How much is your glucose level?


50211272

What's your smoking history

Never smoked

Predict

You have no chance of getting stroke 😊



Made with Streamlit

← Manage app

×

Select your gender

☒ Male
 ☐ Female

Do you have hypertension?

☒ Yes
 ☐ No

Do you have heart disease?

☒ Yes
 ☐ No

Have you ever been married?

☒ Yes
 ☐ No

Heart ♥ Stroke Prediction App

Select your residence type

Urban

How old are you?

8

79

100

Select the type of work you do.

Govt job

How much is your bmi?

10

82

100

How much is your glucose level?

50

243


272

What's your smoking history

Smokes

Predict

You are at risk of getting stroke🙄



Made with Streamlit

← Manage app

← Manage app

6. CONCLUSION

As heart diseases and strokes are increasing rapidly worldwide and causing deaths, it becomes necessary to develop an efficient system that would predict the heart stroke effectively beforehand so that immediate medical attention can be given. In the proposed method, the most effective algorithm for stroke prediction was obtained after a comparative analysis of the accuracy scores of various models. The most effective one was Random Forest, with an accuracy score of 90%.


```
(1, 0, 0, 0),
(1, 0, 0, 0),
(1, 0, 0, 0)]

In [271]: df['avg_glucose_level'].unique()

Out[271]: array(['gl_over 75%', 'gl_50% to 75%', 'gl_under 25%', 'gl_2% to 50%'],
      dtype=object)

mapping
• columns: ['gl_26% to 50%', 'gl_50% to 75%', 'gl_lover 75%', 'gl_under 25%']

In [272]: pd.DataFrame(one_hot_encoded, columns=['gl_26% to 50%', 'gl_50% to 75%', 'gl_lover 75%', 'gl_under 25%'])

Out[272]:
   gl_26% to 50%  gl_50% to 75%  gl_lover 75%  gl_under 25%
0              0              0              1              0
1              0              0              1              0
2              0              1              0              0
3              0              0              1              0
4              0              0              1              0
...          ...          ...          ...          ...
5104           1              0              0              0
5105           0              0              1              0
5106           1              0              0              0
5107           0              0              1              0
5108           1              0              0              0
5109 rows x 4 columns

In [273]: agl_df = pd.DataFrame(one_hot_encoded, columns=['gl_26% to 50%', 'gl_50% to 75%', 'gl_

In [274]: df['bmi']

Out[274]:
0      bmi_over 75%
1      bmi_50% to 75%
2      bmi_50% to 75%
3      bmi_lover 75%
4      bmi_26% to 50%
...
5105      bmi_50% to 75%
5106      bmi_lover 75%
5107      bmi_50% to 75%
5108      bmi_26% to 50%
5109      bmi_26% to 50%
Name: bmi, Length: 5109, dtype: object

In [275]: one_hot_encoded = lb.fit_transform(df['bmi'])
one_hot_encoded

Out[275]: array([[0, 0, 1, 0],
      [0, 1, 0, 0],
      [0, 1, 0, 0],
      [0, 1, 0, 0],
      [0, 1, 0, 0],
      [1, 0, 0, 0]])

In [276]: df['bmi'].unique()

Out[276]: array(['bmi_lover 75%', 'bmi_50% to 75%', 'bmi_26% to 50%',
      'bmi_under 25%'], dtype=object)

mapping
• columns: ['bmi_26% to 50%', 'bmi_50% to 75%', 'bmi_lover 75%', 'bmi_under 25%']

In [277]: pd.DataFrame(one_hot_encoded, columns=['bmi_26% to 50%', 'bmi_50% to 75%', 'bmi_lover

Out[277]:
   bmi_26% to 50%  bmi_50% to 75%  bmi_lover 75%  bmi_under 25%
0              0              0              1              0
1              0              1              0              0
2              0              0              1              0
3              1              0              0              0
4              0              0              0              0
...          ...          ...          ...          ...
5104           0              1              0              0
5105           0              0              1              0
5106           0              1              0              0
5107           1              0              0              0
5108           1              0              0              0
5109 rows x 4 columns

In [278]: bmi_df = pd.DataFrame(one_hot_encoded, columns=['bmi_26% to 50%', 'bmi_50% to 75%', 'b

In [279]: df['smoking_status']

Out[279]:
0      formerly smoked
1      never smoked
2      never smoked
3      smokes
4      never smoked
...
5105      never smoked
5106      never smoked
5107      never smoked
5108      formerly smoked
5109      Unknown
Name: smoking_status, Length: 5109, dtype: object

In [280]: df['smoking_status'].unique()

Out[280]: array(['formerly smoked', 'never smoked', 'smokes', 'Unknown'],
      dtype=object)

In [281]: lb.fit_transform(df['smoking_status'])

Out[281]: array([[0, 1, 0, 0],
      [0, 0, 1, 0],
      [0, 0, 1, 0],
      [0, 0, 1, 0],
      [0, 1, 0, 0],
      [0, 1, 0, 0]])

In [282]: one_hot_encoded = lb.fit_transform(df['smoking_status'])

mapping
• columns: ['Unknown', 'formerly smoked', 'never smoked', 'smokes']

In [283]: pd.DataFrame(one_hot_encoded, columns=['Unknown', 'formerly smoked', 'never smoked', 'sm

Out[283]:
   Unknown  formerly smoked  never smoked  smokes
0         0              1              0         0
1         0              0              1         0
2         0              0              1         0
3         1              0              0         0
4         0              0              1         0
...      ...          ...          ...      ...
5104      0              0              1         0
5105      0              0              1         0
5106      0              1              0         0
5107      1              0              0         0
5108      1              0              0         0
5109 rows x 4 columns

In [284]: smoked_df = pd.DataFrame(one_hot_encoded, columns=['Unknown', 'formerly smoked', 'never

In [285]: cat_cols.append('age')

In [286]: cat_cols

Out[286]: ['age', 'work_type', 'avg_glucose_level', 'bmi', 'smoking_status', 'age']

finally merge!

In [287]: df

Out[287]:
   id  gender  age  hypertension  heart_disease  ever_married  work_type  Residence_type  avg_glu
0    9046      1  age_over
         60              0              1              1  Private              1  g
1    51676      0  age_over
         60              0              0              1  Self-employed              0  g
2    31112      1  age_over
         60              0              1              1  Private              0  gL5
3    60182      0  age_40
         to 59              0              0              1  Private              1  g
4    16655      0  age_over
         60              1              0              1  Self-employed              0  g
...      ...      ...          ...          ...          ...      ...      ...
5105    18234      0  age_over
         60              1              0              1  Private              1  gL2
5106    44873      0  age_over
         60              0              0              1  Self-employed              1  g
5107    19723      0  age_20
         to 39              0              0              1  Self-employed              0  gL2
5108    37544      1  age_40
         to 59              0              0              1  Private              0  g
5109    44679      0  age_40
         to 59              0              0              1  Govt_job              1  gL2
5109 rows x 12 columns

In [288]: df.drop(cat_cols,axis=1,inplace=True)
df

Out[288]:
   id  gender  hypertension  heart_disease  ever_married  Residence_type  stroke
0    9046      1              0              1              1              1              1
1    51676      0              0              0              1              0              1
2    31112      1              0              1              1              0              1
3    60182      0              0              0              1              1              1
4    16655      0              1              0              1              0              1
...      ...      ...          ...          ...          ...      ...
5105    18234      0              1              0              1              1              0
5106    44873      0              0              0              1              1              0
5107    19723      0              0              0              1              0              0
5108    37544      1              0              0              1              0              0
5109    44679      0              0              0              1              1              0
5109 rows x 7 columns

In [289]: # delete meaningless column (for machie learning)
df.drop('id', axis=1, inplace=True)

In [290]: df

Out[290]:
   gender  hypertension  heart_disease  ever_married  Residence_type  stroke
0         1              0              1              1              1              1
1         0              0              0              1              0              1
2         1              0              1              1              0              1
3         0              0              0              1              1              1
4         0              1              0              1              0              1
...      ...          ...          ...          ...      ...
5105      0              1              0              1              1              0
5106      0              0              0              1              1              0
5107      0              0              0              1              0              0
5108      1              0              0              1              0              0
5109      0              0              0              1              1              0
5109 rows x 6 columns

In [291]: cat_cols

Out[291]: ['age', 'work_type', 'avg_glucose_level', 'bmi', 'smoking_status', 'age']

In [292]: df.reset_index().iloc[:,1:]

Out[292]:
   gender  hypertension  heart_disease  ever_married  Residence_type  stroke
0         1              0              1              1              1              1
1         0              0              0              1              0              1
2         1              0              1              1              0              1
3         0              0              0              1              1              1
4         0              1              0              1              0              1
...      ...          ...          ...          ...      ...
5104      0              0              0              1              1              0
5105      0              0              0              1              1              0
5106      0              0              0              1              0              0
5107      1              0              0              1              0              0
5108      1              0              0              1              0              0
5109      0              0              0              1              1              0
5109 rows x 6 columns

In [293]: df = df.reset_index().iloc[:,1:]

In [294]: df

Out[294]:
   gender  hypertension  heart_disease  ever_married  Residence_type  stroke
0         1              0              1              1              1              1
1         0              0              0              1              0              1
2         1              0              1              1              0              1
3         0              0              0              1              1              1
4         0              1              0              1              0              1
...      ...          ...          ...          ...      ...
5104      0              1              0              1              1              0
5105      0              0              0              1              1              0
5106      0              0              0              1              0              0
5107      1              0              0              1              0              0
5108      0              0              0              1              0              0
5109 rows x 6 columns

In [295]: pd.concat([df,age_df,agl_df,work_type_df,bmi_df,smoked_df], axis=1)

Out[295]:
   gender  hypertension  heart_disease  ever_married  Residence_type  stroke  age_20
to 39  age_40
to 59  age_over
to 60
0         1              0              1              1              1              1              0              0              1
1         0              0              0              1              0              1              0              0              1
2         1              0              1              1              0              1              0              0              1
3         0              0              0              1              1              1              0              0              1
4         0              1              0              1              0              1              0              0              1
...      ...          ...          ...          ...      ...      ...
5104      0              1              0              1              1              0              0              0              1
5105      0              0              0              1              1              0              0              0              1
5106      0              0              0              1              0              0              0              0              1
5107      1              0              0              1              0              0              0              1              0
5108      0              0              0              1              1              0              0              0              1
5109 rows x 27 columns

In [296]: pd.concat([df,age_df,agl_df,work_type_df,bmi_df,smoked_df], axis=1).isnull().sum()

Out[296]:
gender              0
hypertension        0
heart_disease       0
ever_married        0
Residence_type      0
stroke              0
age_20 to 39        0
age_40 to 59        0
age_over 60         0
age under 19        0
gl_2% to 50%        0
gl_50% to 75%       0
gl_lover 75%        0
gl_under 25%        0
Govt_job            0
Never_worked        0
Private             0
Self-employed       0
Children            0
bmi_26% to 50%      0
bmi_50% to 75%      0
bmi_lover 75%       0
bmi_under 25%       0
Unknown            0
formerly smoked     0
never smoked        0
smokes              0
dtype: int64

In [297]: final_df = pd.concat([df,age_df,agl_df,work_type_df,bmi_df,smoked_df], axis=1)

Step 3. Train, Test set split & Upsampling

In [298]: from sklearn.model_selection import train_test_split

In [299]: X = final_df.drop('stroke', axis=1)
y = final_df['stroke']

In [300]: X.shape, y.shape

Out[300]: ((5109, 26), (5109,))

In [301]: X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state=11)

In [302]: X_train.shape, y_train.shape

Out[302]: ((4087, 26), (4087,))

In [303]: X_test.shape, y_test.shape

Out[303]: ((1022, 26), (1022,))

In [304]: y_train.value_counts()

Out[304]:
0    3893
1     194
Name: stroke, dtype: int64

In [305]: y_test.value_counts()

Out[305]:
0     967
1      55
Name: stroke, dtype: int64

In [306]: plt.bar(x = y_train.value_counts().index, height = y_train.value_counts().values);



It's imbalanced data, we have to resample '1' in training set-y

In [307]: from sklearn.utils import resample

In [308]: train_df = pd.concat([X_train,y_train], axis=1)

In [309]: train_df

Out[309]:
   gender  hypertension  heart_disease  ever_married  Residence_type  age_20
to 39  age_40
to 59  age_over
to 60
1977      1              0              0              1              1              1              0              0
3378      0              0              0              1              1              0              1              0
924       0              0              0              1              1              0              0              1
4104      0              0              0              0              1              0              1              0
3878      0              0              0              1              0              1              0              0
...      ...          ...          ...          ...      ...      ...
681       0              0              0              0              0              0              0              0
4182      0              0              0              1              0              0              0              1
4820      0              0              0              1              0              1              0              0
2004      1              0              0              1              0              0              1              0
3924      0              0              0              1              0              0              0              1
4087 rows x 27 columns

In [310]: train_0 = train_df[train_df['stroke']==0]
train_1 = train_df[train_df['stroke']==1]

In [311]: train_0.shape, train_1.shape

Out[311]: ((3893, 27), (194, 27))

In [312]: upsampled_train_1 = resample(train_1,
                                  replace=True,
                                  n_samples=3893,
                                  random_state=123)

Out[312]: (3893, 27)

In [313]: upsampled_train = pd.concat([train_0, upsampled_train_1])

In [314]: upsampled_train['stroke'].value_counts()

Out[314]:
1    3893
0    3893
Name: stroke, dtype: int64

In [315]: X_train = upsampled_train.drop('stroke',axis=1)
y_train = upsampled_train['stroke']

Step 4. Modeling & Prediction

1. RandomForestClassifier

In [316]: from sklearn.ensemble import RandomForestClassifier

In [317]: rfc = RandomForestClassifier(n_estimators=100)
rfc.fit(X_train, y_train)
print('Train Accuracy : {:.2f}'.format(rfc.score(X_train, y_train)))
print('Test Accuracy : {:.2f}'.format(rfc.score(X_test, y_test)))

Train Accuracy : 0.97
Test Accuracy : 0.90

1. GradientBoostingClassifier

In [318]: best_params = {}
score = 0
for i in range(1,8):
    for j in [50, 100, 150, 200, 250, 300, 350, 400]:
        gbc = GradientBoostingClassifier(max_depth=i,
                                         n_estimators=j)
        gbc.fit(X_train, y_train)
        print('max_depth : {}'.format(i))
        print('n_estimators : {}'.format(j))
        rfc_score = rfc.score(X_train, y_train)
        print('Train Score : {}'.format(gbc.score(X_train, y_train)))
        print('Test Score : {}'.format(gbc.score(X_test, y_test)))
        print('-----')
        if gbc.score(X_test, y_test) > score:
            score = gbc.score(X_test, y_test)
            best_params['max_depth'] = i
            best_params['n_estimators'] = j
        best_params['n_estimators'] = j

max_depth : 1
n_estimators : 50
Train Score : 0.74939357821731313
Train Score : 0.7426614981409002
-----
max_depth : 1
n_estimators : 100
Train Score : 0.7528898022909933
Train Score : 0.7172211350293543
-----
max_depth : 1
n_estimators : 150
Train Score : 0.7549447726689929
Train Score : 0.7084148727984344
-----
max_depth : 1
n_estimators : 200
Train Score : 0.7537888517852556
Train Score : 0.7074363992172211
-----
max_depth : 1
n_estimators : 250
Train Score : 0.7549447726689929
Train Score : 0.7045067514677103
-----
max_depth : 1
n_estimators : 300
Train Score : 0.7566144361674801
Train Score : 0.7035225048923679
-----
max_depth : 1
n_estimators : 350
Train Score : 0.7562375649600052
Train Score : 0.702544031311546
-----
max_depth : 1
n_estimators : 400
Train Score : 0.7564860005137426
Train Score : 0.702544031311546
-----
max_depth : 2
n_estimators : 50
Train Score : 0.7634715259155664
Train Score : 0.71526481878669276
-----
max_depth : 2
n_estimators : 100
Train Score : 0.77125016006935525
Train Score : 0.6966731898238747
-----
max_depth : 2
n_estimators : 150
Train Score : 0.7634715259155664
Train Score : 0.6956947162426614
-----
max_depth : 2
n_estimators : 200
Train Score : 0.778147957873105
Train Score : 0.702544031311546
-----
max_depth : 2
n_estimators : 300
Train Score : 0.7954020035961983
Train Score : 0.7084148727984344
-----
max_depth : 2
n_estimators : 350
Train Score : 0.8006678653994349
Train Score : 0.7084148727984344
-----
max_depth : 2
n_estimators : 400
Train Score : 0.8115948959671205
Train Score : 0.7093933463796478
-----
max_depth : 3
n_estimators : 50
Train Score : 0.7819162599537631
Train Score : 0.7064579256360078
-----
max_depth : 3
n_estimators : 100
Train Score : 0.8198232725404573
Train Score : 0.7084148727984344
-----
max_depth : 3
n_estimators : 150
Train Score : 0.8217313126123812
Train Score : 0.7201565557729941
-----
max_depth : 3
n_estimators : 200
Train Score : 0.8493494781659389
Train Score : 0.7416829745596689
-----
max_depth : 3
n_estimators : 300
Train Score : 0.8517852535869501
Train Score : 0.7475538160469667
-----
max_depth : 3
n_estimators : 350
Train Score : 0.8594913948111996
Train Score : 0.7534246573542466
-----
max_depth : 3
n_estimators : 400
Train Score : 0.8688671975340354
Train Score : 0.7622309197651063
-----
max_depth : 4
n_estimators : 50
Train Score : 0.8217313126123812
Train Score : 0.7328767123287672
-----
max_depth : 4
n_estimators : 100
Train Score : 0.8556383253390752
Train Score : 0.7475538160469667
-----
max_depth : 4
n_estimators : 150
Train Score : 0.8835088620601079
Train Score : 0.7661448140900196
-----
max_depth : 4
n_estimators : 200
Train Score : 0.93128737143591061
Train Score : 0.782778649706457
-----
max_depth : 4
n_estimators : 250
Train Score : 0.9072694580015412
Train Score : 0.7945205479452054
-----
max_depth : 4
n_estimators : 300
Train Score : 0.9170305676855895
Train Score : 0.797455968888454
-----
max_depth : 4
n_estimators : 350
Train Score : 0.9235807860262009
Train Score : 0.8072407450409785
-----
max_depth : 4
n_estimators : 400
Train Score : 0.93128737143591069
Train Score : 0.812133072407045
-----
max_depth : 5
n_estimators : 50
Train Score : 0.8751605446571718
Train Score : 0.7720156557727995
-----
max_depth : 5
n_estimators : 100
Train Score : 0.9140765476496275
Train Score : 0.8043052837573386
-----
max_depth : 5
n_estimators : 150
Train Score : 0.9287182121756999
Train Score : 0.8209393346379648
-----
max_depth : 5
n_estimators : 200
Train Score : 0.9407911636270229
Train Score : 0.8228962180003914
-----
max_depth : 5
n_estimators : 250
Train Score : 0.9501669634498588
Train Score : 0.8426615734246576
-----
max_depth : 5
n_estimators : 300
Train Score : 0.9546622142306704
Train Score : 0.8512720156555773
-----
max_depth : 5
n_estimators : 350
Train Score : 0.960955612638069
Train Score : 0.85318595099804305
-----
max_depth : 5
n_estimators : 400
Train Score : 0.9642948882609812
Train Score : 0.860078277886497
-----
max_depth : 6
n_estimators : 50
Train Score : 0.9081685075777036
Train Score : 0.797455968888454
-----
max_depth : 6
n_estimators : 100
Train Score : 0.9277716414076548
Train Score : 0.8933894324853229
-----
max_depth : 6
n_estimators : 150
Train Score : 0.93128737143591069
Train Score : 0.8228962180003914
-----
max_depth : 6
n_estimators : 200
Train Score : 0.9681479578731056
Train Score : 0.87179843444227
-----
max_depth : 6
n_estimators : 250
Train Score : 0.9681479578731056
Train Score : 0.87179843444227
-----
max_depth : 6
n_estimators : 300
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 6
n_estimators : 350
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 6
n_estimators : 400
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 50
Train Score : 0.93128737143591069
Train Score : 0.8228962180003914
-----
max_depth : 7
n_estimators : 100
Train Score : 0.9681479578731056
Train Score : 0.87179843444227
-----
max_depth : 7
n_estimators : 150
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 200
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 250
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 300
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 350
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 400
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 450
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 500
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 550
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 600
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 650
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 700
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 750
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 800
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 850
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 900
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 950
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1000
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1050
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1100
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1150
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1200
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1250
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1300
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1350
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1400
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1450
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1500
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1550
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1600
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1650
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1700
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1750
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1800
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1850
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1900
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 1950
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2000
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2050
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2100
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2150
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2200
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2250
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2300
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2350
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2400
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2450
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2500
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2550
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2600
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2650
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2700
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2750
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2800
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2850
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2900
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 2950
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 3000
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 3050
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 3100
Train Score : 0.97342253300797
Train Score : 0.8933463796477495
-----
max_depth : 7
n_estimators : 3150
Train Score : 0.97342253300797
Train Score : 0.8933463796477
```

```

import pandas as pd
import numpy as np
import pickle
import streamlit as st
from PIL import Image

pickle_in = open('classifier.pkl', 'rb')
classifier = pickle.load(pickle_in)

def welcome():
    return 'welcome all'

def prediction(gender, age1, hypertension, heart_disease, married, work_type1, residence, avg_glucose1, bmi1, smoking_stat1):
    age = [0,0,0,0]
    work_type = [0,0,0,0,0]
    glucose_level = [0,0,0,0]
    bmi = [0,0,0,0]
    smoking = [0,0,0,0]
    if gender=="Male":
        gender=0
    else:
        gender=1

    if hypertension=="Yes":
        hypertension=1
    else:
        hypertension=0

    if heart_disease=="Yes":
        heart_disease=1
    else:
        heart_disease=0

    if married=="Yes":
        married=1
    else:
        married=0

    if residence=="Rural":
        residence=1
    else:
        residence=0

    if age1>=20 and age1<40:
        age[0]=1
    elif age1>=40 and age1<60:
        age[1]=1
    elif age1>=60:
        age[2]=1
    else:
        age[3]=1

    if avg_glucose1<77:
        glucose_level[3]=1
    elif avg_glucose1<91:
        glucose_level[0]=1
    elif avg_glucose1<114:
        glucose_level[1]=1
    else:
        glucose_level[2]=1

    if work_type1=="Govt job":
        work_type[0]=1
    elif work_type1=="Never worked":
        work_type[1]=1
    elif work_type1=="Private":
        work_type[2]=1
    elif work_type1=="Self-employed":
        work_type[3]=1
    else:
        work_type[4]=1

    if bmi1<23:
        bmi[3]=1
    elif bmi1<28:
        bmi[0]=1
    elif bmi1<33:
        bmi[1]=1
    else:
        bmi[2]=1

    if smoking_stat1=="Unknown":
        smoking[0]=1
    elif smoking_stat1=="Formerly smoked":
        smoking[1]=1
    elif smoking_stat1=="Never smoked":
        smoking[2]=1
    else:

```

```

smoking[3]=1

prediction = classifier.predict(

[[gender,hypertension,heart_disease,married,residence,age[0],age[1],age[2],age[3],glucose_level[0],glucose_level[1],glucose_level[2],glucose_level[3],

work_type[0],work_type[1],work_type[2],work_type[3],work_type[4],bmi[0],bmi[1],bmi[2],bmi[3],smoking[0],smoking[1],smoking[2],smoking[3]])]

#print(heart_disease)
return prediction

def main():
    st.title("Heart ❤️ Stroke Prediction App")
    html_temp=""
    ans=0
    st.markdown(html_temp,unsafe_allow_html = True)
    get_Gender = st.sidebar.radio("Select your gender",("Male","Female"))
    get_Hypertension = st.sidebar.radio("Do you have hypertension?",("Yes","No"))
    get_heartDisease = st.sidebar.radio("Do you have heart disease?",("Yes","No"))
    get_married = st.sidebar.radio("Have you ever been married?",("Yes","No"))
    get_residence = st.selectbox("Select your residence type",("Rural","Urban"))
    get_age = st.slider("How old are you?",value=25)
    get_workType = st.selectbox("Select the type of work you do.",("Govt job","Never worked","Private","Self-employed","Children"))
    get_bmi = st.slider("How much is your bmi?",min_value=10,max_value=100,value=65)
    get_glucose = st.slider("How much is your glucose level?",min_value=55,max_value=272,value=50)
    get_smoke = st.selectbox("What's your smoking history",("Unknown","Formerly smoked","Never smoked","Smokes"))

    if st.button("Predict"):

ans=prediction(get_Gender,get_age,get_Hypertension,get_heartDisease,get_married,get_workType,get_residence,get_glucose,get_bmi,get_smoke)[0]
        if ans==0:
            st.success('You have no chance of getting stroke😊')
            st.image('images/happy_heart.jfif')
        else:
            st.success('You are at risk of getting stroke😞')
            st.image('images/damaged_heart.jfif')

if __name__=='__main__':
    main()

```