

Bellevue University

White Paper On
Can we predict Stroke?

By:
Debasish Panda

Introduction:

A stroke, or brain attack, happens when blood flow to our brain is stopped. It is an emergency situation. The brain needs a constant supply of oxygen and nutrients in order to work well. If blood supply is stopped even for a short time, this can cause problems. Brain cells begin to die after just a few minutes without blood or oxygen. When brain cells die, brain function is lost. We may not be able to do things that are controlled by that part of the brain. For example, a stroke may affect our ability to move, speak, eat, think and remember, controlling emotions, controlling bowel and bladder etc. In other words, a stroke can happen to anyone at any time.

- According to CDC, every year more than 795,000 people in United States have a stroke. Every 4 minutes someone dies of stroke.
- Even though stroke is a leading cause of death for Americans, but the risk of having a first stroke is nearly twice as high for blacks as for whites, and blacks have the highest rate of death due to stroke.
- Recurrent stroke is frequent; about 25 percent of people who recover from their first stroke will have another stroke within 5 years.

Since stroke is something that can happen to anyone at any age, if we can use some of the prediction modelling techniques to predict whether a given person is likely to have a stroke or not based on the risk factors, we discussed then we can save a lot more lives.

Scope:

The main objective of this paper is to design a robust system that works efficiently and will be able to predict the possibility of Stroke accurately. This paper uses the dataset available in Kaggle with 5110 observations.

Link to dataset: <https://www.kaggle.com/datasets/fedesoriano/stroke-prediction-dataset>

Here are the attributes of the dataset explained:

- 1) id: unique identifier
 - 2) gender: "Male", "Female" or "Other"
 - 3) age: age of the patient
 - 4) hypertension: 0 if the patient doesn't have hypertension, 1 if the patient has hypertension
 - 5) heart_disease: 0 if the patient doesn't have any heart diseases, 1 if the patient has a heart disease
 - 6) ever_married: "No" or "Yes"
 - 7) work_type: "children", "Govt_job", "Never_worked", "Private" or "Self-employed"
 - 8) Residence_type: "Rural" or "Urban"
 - 9) avg_glucose_level: average glucose level in blood
 - 10) bmi: body mass index
 - 11) smoking_status: "formerly smoked", "never smoked", "smokes" or "Unknown"*
 - 12) stroke: 1 if the patient had a stroke or 0 if not
- Below is the link to the actual data set.

Methodology:

Technical approach

This course project on predictive analytics on heart failure will follow the CRISP-DM model for data understanding, modeling, and evaluation. I will use the Python programming language to load, explore, and model the prediction system along with necessary machine learning libraries for Python.

Data Analysis

Initially, exploratory data analysis was performed to understand the data better and compute some major descriptive statistics. I produced some illustrative visualizations representing the dataset as well. Then, various machine learning algorithms were employed to see the best- resulting model for the prediction model. Mainly, SVM, Logistic Regression, Decision Tree, Random Forest, and KNN are the potential algorithms used.

The collected dataset is analyzed using Python programming language, and Jupyter Notebook is used as a scripting interface. Several python libraries are used for the various purposes of exploratory data analysis, visualization, and machine learning model creation. Some of the notable ones are:

- Numpy
- Pandas
- Matplotlib
- Seaborn
- Sklearn

After loading the dataset into the pandas data frame, different measures about the dataset were generated, such as the shape of the dataset, data types, missing values, duplicated observations, etc. Then Exploratory data analysis was performed using visualization tools. Then the critical step of the project, modeling, was performed. Finally, data scaling, numerical and categorical data, and encoding were accomplished.

Data was trained with several classification algorithms as the data processing was complete. Trained models were attempted to improve by hyper-parameter tuning. Data was first trained with Logistic Regression, SVM, GaussianNB, BernoulliNB, Decision tree, Random Forest classifier- Nearest neighbor. Fine-tuning of each of those models was performed and evaluated in each step with testing data. The model evaluation mainly used accuracy, F1 score, and whole classification report. In the case of the KNN classifier, various k values were tested, and the best one was selected to train the model.

Results

Exploratory Data Analysis

Dataset was imported in Python's Jupyter notebook. Initially, the dataset was imported as a panda's data frame. Here's what the dataset looked like.

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type	avg_glucose_level	bmi	smoking_status	stroke
0	9046	Male	67.0	0	1	Yes	Private	Urban	228.69	36.6	formerly smoked	1
1	51676	Female	61.0	0	0	Yes	Self-employed	Rural	202.21	NaN	never smoked	1
2	31112	Male	80.0	0	1	Yes	Private	Rural	105.92	32.5	never smoked	1
3	60182	Female	49.0	0	0	Yes	Private	Urban	171.23	34.4	smokes	1
4	1665	Female	79.0	1	0	Yes	Self-employed	Rural	174.12	24.0	never smoked	1

Figure 1 : Dataset preview

Dataset has 12 columns and 5110 rows of observations. As soon as the data was loaded, missing and duplicated values were checked. It was found that there were 201 rows missing values in bmi. I have imputed the missing values with the median value of the column.

There is one row with gender value as "other" (not Male and Female). I have removed that row from the data set as a cleanup.

Below is a histogram of numerical variables in the feature set.

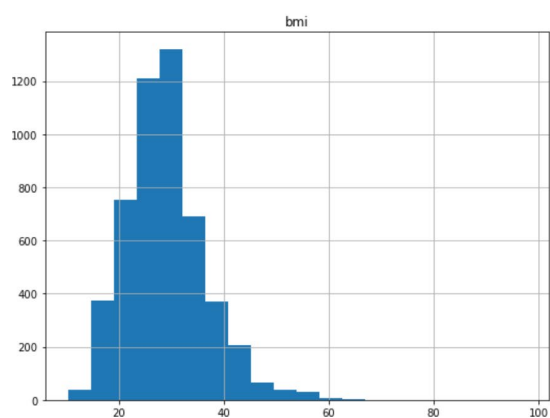
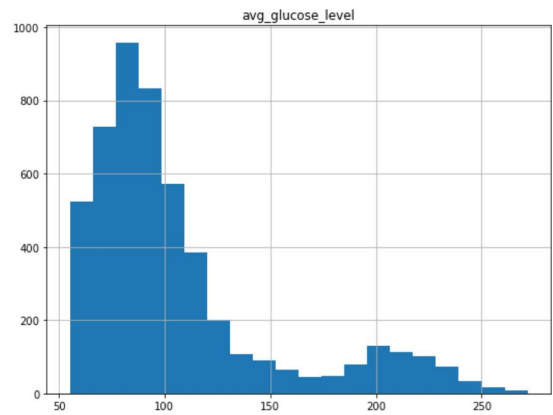
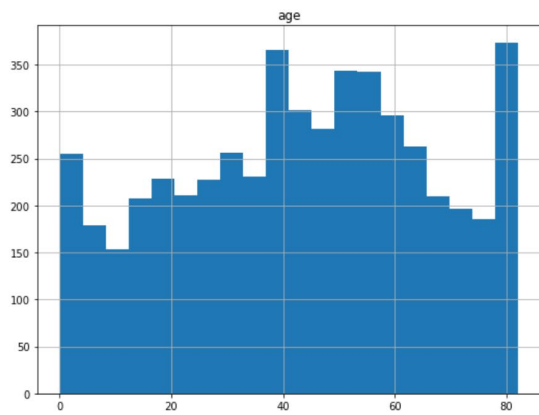


Figure 2: Histograms of Numerical variables

Figure 2 shows that we have a reasonable variation in the age in the input data set. Also, none of the numerical variables have outliers to clean. For example, glucose value and BMI are centered around 90 and 30 weights, respectively. This data behavior is expected as most incoming patients tend to be centered around this mode value as slightly borderline glucose (100 max) and Overweight/obese bmi range of 30.

Here is a heat map correlation of categorical variables against Stroke.

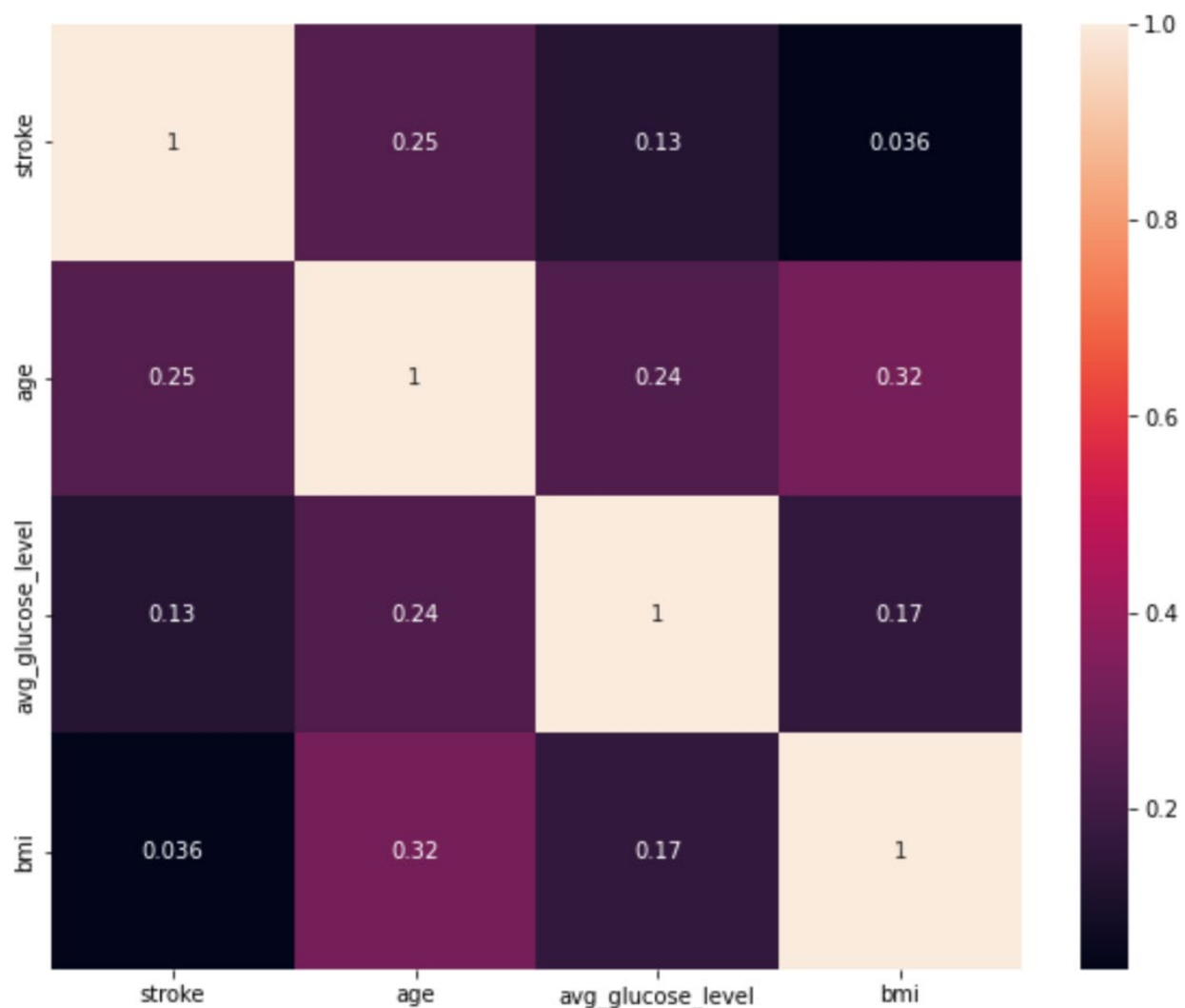


Figure 3: Heat Correlation against numerical features

Figure 3 shows a strong correlation with some of the categorical variables like age against stroke prediction.

Here is Figure 4, outlining the correlation heat map of categorical features against Stroke.

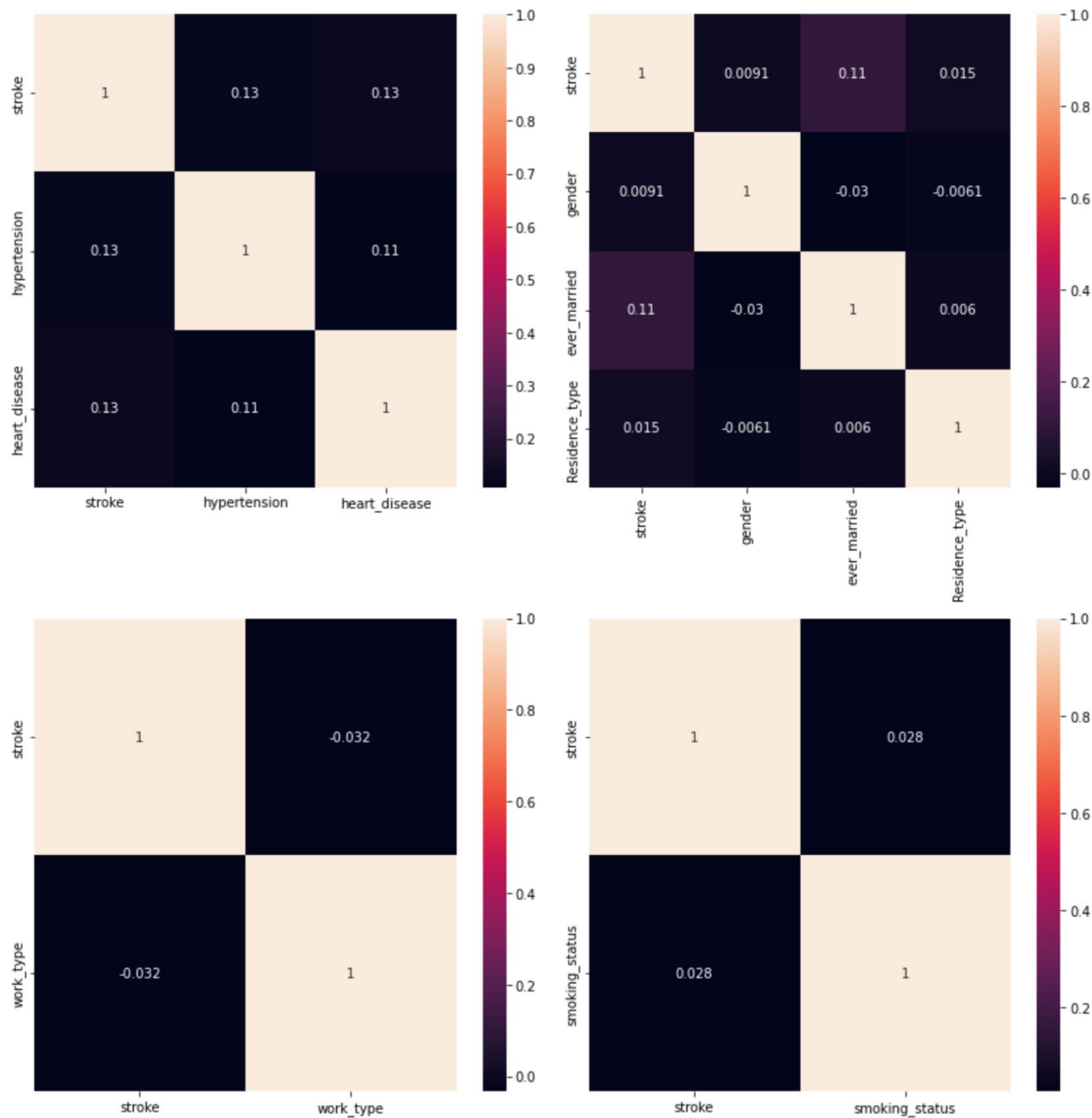


Figure 4: Correlation heatmap for categorical variables

Figure 4 shows a correlation of categorical values against Stroke. Some slightly significant correlations with prior heart disease hypertension if the patient is married. Very first look, it shows as if there is no significant correlation with residence type, work type, or smoking status of the patient. But a lot of anecdotal evidence states that stress in life adds risk factors to Stroke.

Also, based on several medical research outcomes, smoking is one of the severe risk factors for heart diseases and hypertension. As these two factors seem to have some significant correlation with Stroke, I would like to retain smoking status and an input feature to my model. This could indicate an impending cascading event of heart disease or Stroke.

After the exploratory data analysis, the focus shifted to model development. The first model built was using the KNN classifier with a grid search using 1-10 k value search space. The model showed significant accuracy and corresponding scores, with k=8 as best performing. However, a closer look at the confusion matrix proved that the model is not identifying true positives in the target test set. Here is a look at the model outcome.

```
Accuracy: 0.9515
precision: 0.9538
recall: 0.9515
f1 Score: 0.9286
Confusion Matrix for Prediction:

array([[1215,    0],
       [  62,    1]])
```

Figure5: KNN model performance metrics before Class balancing

This seemed like a result of imbalanced classes in the train data set for my model. I have rebalanced by training data set using SMOTE as a next step. Here is a look at before and after for my train data set,

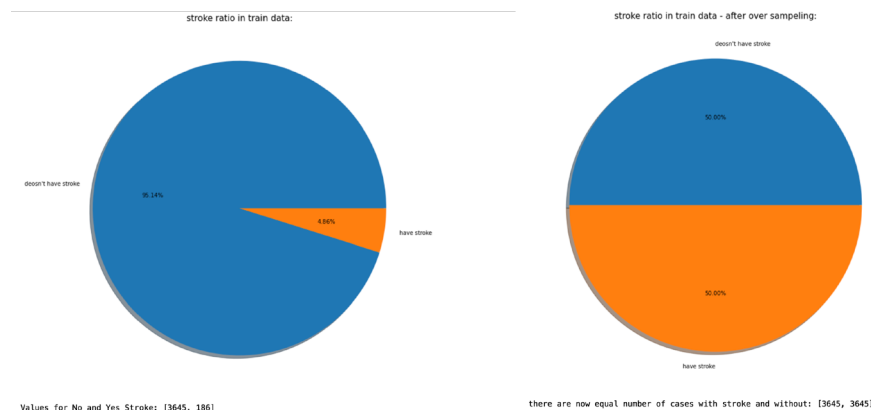


Figure 6: Train data set before vs. After Class balance using SMOTE

After class rebalances, I have trained various models to see the performance of the data set like Logistics regression, SVM, GassianNB, BernoulliNB, KNN classification (without grid search), Decision tree, and Random search. All the models performed with varying degrees of accuracy from 0.56 to 0.89 with different degrees of success in identifying true negatives and positives. Out of all the models, K-nearest neighbor balanced measurement scores and identified true positives and negatives.

After that, the K-nearest neighbor classifier was employed again with a grid search of 1-10. This time, the model performed better by pushing accuracy from 0.81 to 0.89 and keeping true positive and negative identification compared to my original trained KNN model before class rebalancing.

```
Accuracy: 0.8889
precision: 0.9118
recall: 0.8889
f1 Score: 0.8999
Confusion Matrix for Prediction:

array([[1127,   88],
       [  54,    9]])
```

Figure 7: KNN model performance metrics after Class balancing

Discussion

Data understanding

This project aims to create a predictive model and analyze the data of people having a stroke. It is essential to note this dataset itself is not very large and has an Imbalance in target class distribution. But it was something of my interest and the best data I could find out there for free. For a person to develop a disease, hundreds of factors play the role. By no means I could create a model that can predict Stroke based on a few parameters. However, it provided some insights into Stroke disease itself and the process of developing a predictive model.

I should be thankful that the data I acquired was clean with no missing or duplicate values, which saved me a lot of time and effort. I have used mainly the matplotlib and seaborn as the visual tools. They've produced some good visuals in the notebook.

I chose pie charts over other plots for categorical visuals because it's easier to see the distribution and symmetry at a glance. For the numerical variables, though, I generated a set of histograms for each variable.

Data preprocessing

Dealing with missing entries in the bmi column was achieved by imputing the data with a median value, as it's a numerical attribute to avoid skewing the model training. However, my significant effort was to rebalance the imbalanced classes in the train data set. I achieved this by oversampling underrepresented classes by using SMOTE package.

Feature engineering

During the feature engineering, categorical features were encoded using the one-hot encoding method of Pandas library. This enabled the conversion of all the categorical variables into numbers. After that, all the numbers were scaled using standard scalar to transform the data into a scaled attribute.

Model development

Model development was the most critical and insightful part of the project as I had to go through a bunch of classifiers and experiment with them. Unfortunately, some of the classifiers performed so poorly that I did not include that in my notebook.

Even though I knew the data was imbalanced in the first attempt, I have tried a KNN classifier training using grid search hyperparameter tuning for `n_neighbors` to see how the model would perform. This resulted in a higher accuracy of 95% but looking at the confusion matrix of results showed that the model failed to identify any true positive outcomes. This reinforced my learning of the importance of data and review of train data fed into a model.

I had them rebalance the train data set using SMOTE to oversample underrepresented true positive classes. I used this technique as I could not find a free data set that could augment and rebalance my dataset.

Once the data set was rebalanced, I have trained various models to see the performance of the data set like Logistics regression, SVM, GaussianNB, BernoulliNB, KNN classification (without grid search), Decision tree, and Random search. These models have been performed with varying accuracy and true positive and negative identification. Out of all these models, KNN was the most well-balanced.

To boost the model further, I have trained a KNN algorithm with hyperparameter tuning using `n_neighbors`, which resulted in better results of close to 90% accuracy and better true positive and negative identification with precision/recall and f1 scores all-around 90% mark.

Conclusion

Best the predictive model found out to be a KNN classifier with hyperparameter tuning, which has the highest classification accuracy of 90%. I must admit that predicting a stroke is a very sensitive undertaking for anyone. I am confident that this model is the best model for the given limited features as in this study. But it is not the best for predicting a stroke in general because the model doesn't consider other important risk factors like exercise, alcohol consumption, nutrition, etc. Although this project was academic, it gave me practical perspectives on data mining and predictive modeling. Some of the challenging parts of the project were

preprocessing the dataset, dealing with an imbalance in the dataset, fine-tuning the various classification algorithms.

Recommendation

This model can predict stroke disease based on 11 features included in this study with 90 percent accuracy. Still, it can be improved further if more data and features have all other important risk factors of Stroke and a more balanced dataset with higher volume for train data. Therefore, I would recommend any stakeholder use this model with caution.

Questions:

1. Why do we need to predict the Risk of Stroke?

Ans: *Without the blood supply, the brain cells gradually die, and disability occurs depending on the area of the brain affected. Early recognition of symptoms can significantly carry valuable information for the prediction of stroke and promoting a healthy life.*

2. How does this project help in predicting Stroke?

Ans: *This model used KNN classification model to predict the probability of stroke.*

3. Does this disease affect more people of color? If yes, how do you know?

Ans: *Out of the 12 columns present in the dataset there is no data specific to identify race of an individual. Hence, this model does not predict anything specific to any person of color.*

4. What parameters are used in the model?

Ans: *Accuracy, Precision, recall, f1 score, confusion matrix etc are the parameters used in this model.*

5. What are the different models that are built as part of this project?

Ans: *Logistic regression, SVM, KNN, GaussianNB, BernoulliNB, Decision tree, Random forest are the models that are built as part of this project.*

6. Are there any new features built from the existing features?

Ans: *No.*

7. Can adding any other features affect the model accuracy?

Ans: Yes, I have retrained the model using balanced data set which resulted in a little reduction of accuracy but improved the identification of the true positives.

8. is there any additional data to supplement the dataset used?

Ans: No.

9. Why is accuracy a poor model for imbalanced data sets?

Ans: KNN model has resulted in very high accuracy/precision/recall and f1 scores - all of which are in 90%. This high accuracy could be a result of imbalanced dataset (95% negative outcomes, and 5% positive outcomes of stroke)

10. Can this methodology be used in other areas of medical sciences?

Ans: This methodology can be used in other areas of medical sciences.

References:

1. NHLBI. (March24, 2022). What is a Stroke? National Heart, Lung and Bone Institute. Retrieved December 12, 2022, from <https://www.nhlbi.nih.gov/health/stroke>
2. Chris Albon. (2018). *Machine Learning with Python Cookbook Practical Solutions from Preprocessing to Deep Learning*.
3. Stroke Prediction Dataset | Kaggle.
<https://www.kaggle.com/fedesoriano/stroke-prediction-dataset>
4. Stroke background information by CDC from <https://www.cdc.gov/stroke/about.htm>
5. Raj, S. (2020). *How to Evaluate the Performance of Your Machine Learning Model - KDnuggets*. <https://www.kdnuggets.com/2020/09/performance-machine-learning-model.html>
6. scikit-learn. (2021). *API Reference — scikit-learn 1.0.1 documentation*.
<https://scikit-learn.org/stable/modules/classes.html#module-sklearn.metrics>
 - a. SMOTE - <https://machinelearningmastery.com/smote-oversampling-for-imbalanced-classification/>
 - b. Siegel, E. (2013). *Predictive analytics: The power to predict who will click*,

buy, lie, or die. John Wiley & Sons.