# Predicting Survival chances for stroke patients 80 years and above

By: Bhavya Balasubramanya
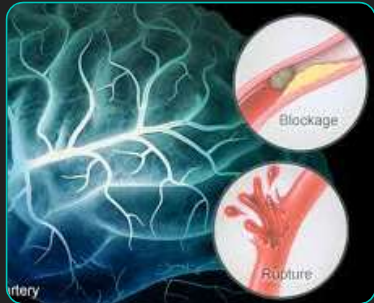
May 2019

For complete project with code, please visit my repository
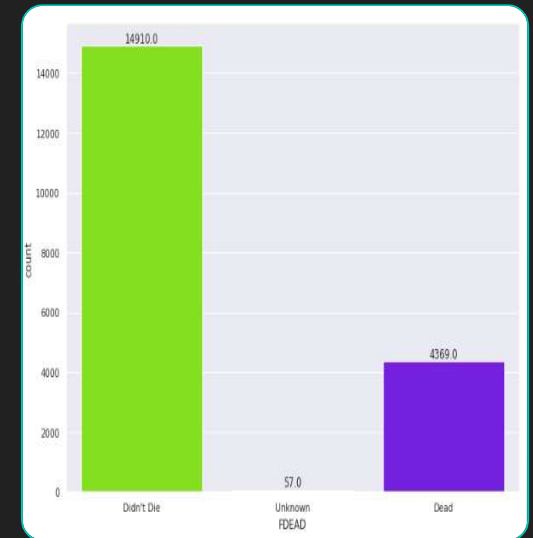
# Outline

# Introduction







According to WHO, 15 million people suffer from stroke every year. Of these 5 million die and other 5 million are permanently disabled. In, U.S., stroke is the 3rd leading cause of death. More than 140,000 people die due to stroke and it is also a leading cause of serious and long-disability in US. On an average someone in US has a stroke every 40 seconds. Risk of having stroke more than doubles each decade after age of 55. It can occur at any age and nearly 1/4th occur in people under age of 65. Risk factors increase with preconditions like diabetes, high blood pressure, high cholesterol, Atrial fibrillation. However, out of these, high blood pressure is the most common factor for stroke. In developed countries, incidence of stroke is declining, largely due to efforts to lower blood pressure and reduced smoking. However, overall rate of stroke remains high due to aging of the population.

# Understanding the dataset

- This dataset has been taken from the International Stroke trial database (version 2) which was one of the biggest randomized trials in actual stroke.

- The data was collected worldwide in 43 different countries.

- There were a total of 19,432 records and 112 attributes.

- The data was recorded from 1991 to 1996.

- Data from 1991 to 1993 was considered pilot phase. We haven't considered this phase in our capstone since there was lot of missing data during this period. So for future reference our data is only from the period March 1993 – April 1996

- Target variable chosen for modelling purposes is 'FDEAD' which indicates if the person is alive or not at sixth month follow up

- The initial distribution of target variable is as shown in Fig A.

- Though there appears to be a bias in data, we will be using this data by removing just the unknowns and not really balancing the data as our predictions actually compares survival rate of patients aged 80 and above to those below.

# Data Cleaning

## 10
### Remove Unwanted columns

There were a lot of columns in our data which was not necessary. Since this data was collected worldwide, the dates didn't make sense as it was not uniform and didn't tell anything important

## 19
### Remove columns and rows that don't have much data

There were some columns for comments which didn't have significant data to use imputation, so they were removed. Also most of the data was missing during pilot phase, so pilot phase rows were removed

## 53
### Convert categorical data to indicator variables

We used dummy data for Stroke types. While doing this, we also replaced most common categorical data like Can't access (C) to 0, No (N) to 1, Yes (Y) to 2 and Unknown (U) to 3

## 1
### Replace unreadable data to readable data

Most of the months in our data had Polish names and English names. Since Polish names are not understandable by the majority, we have converted them to English

# Exploratory Data Analysis

Let us try to explore our data and see how the attributes are related to each other and with our target.

Correlation for Drugs given

Correlation for Randomization Deficit
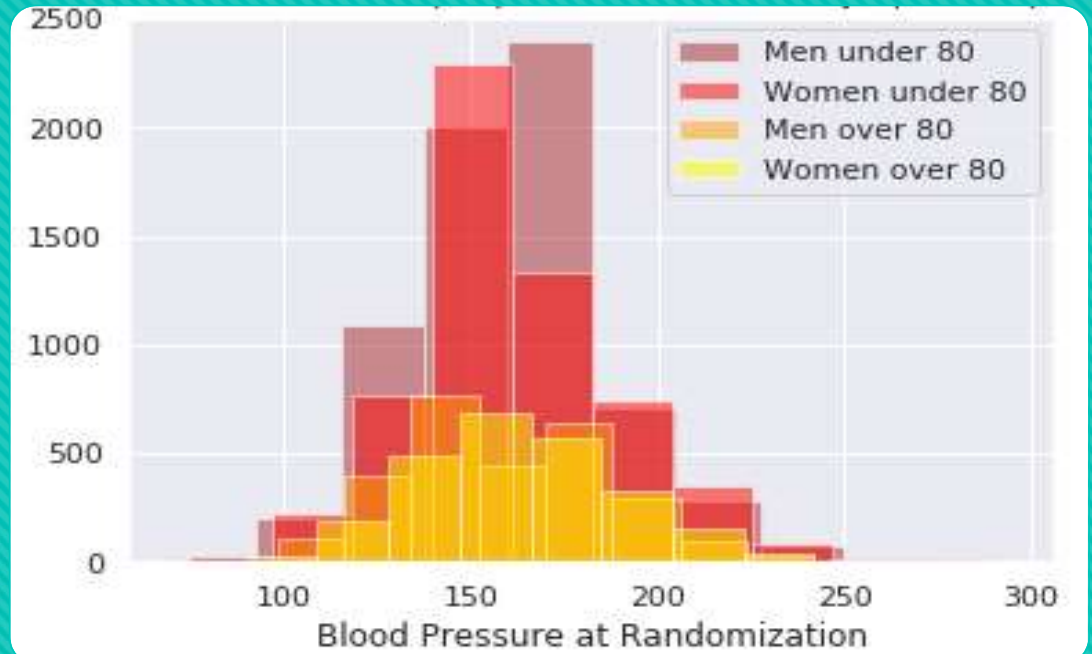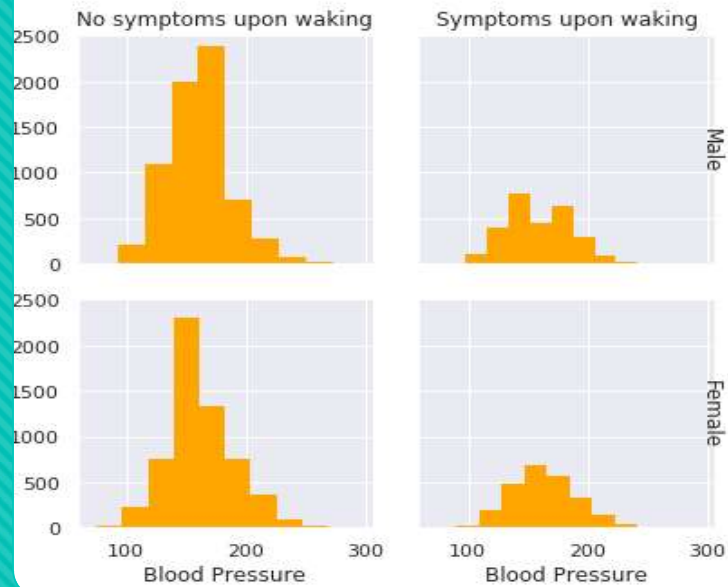
# Correlations

# Correlations Continued..

# Geographical Data
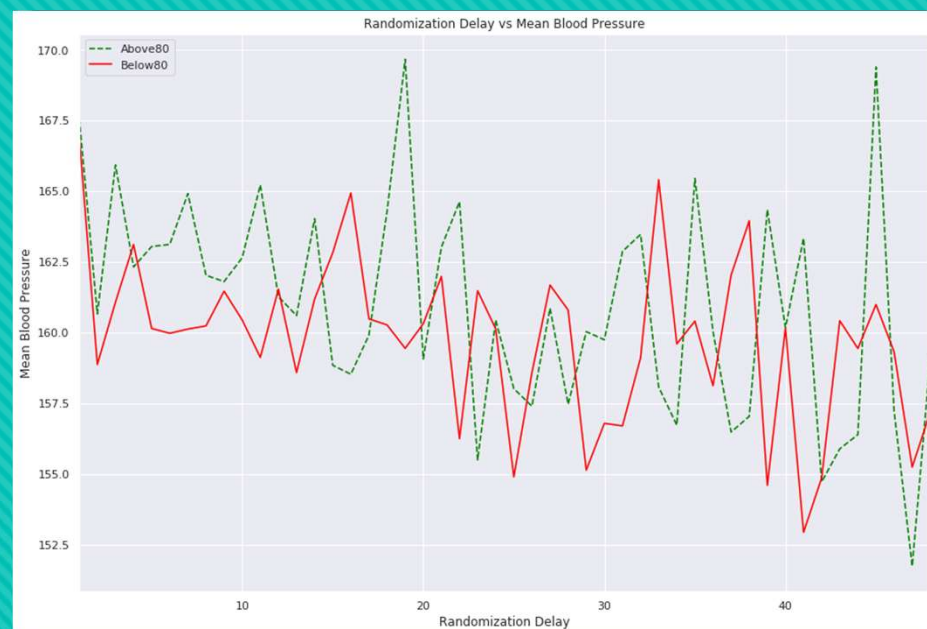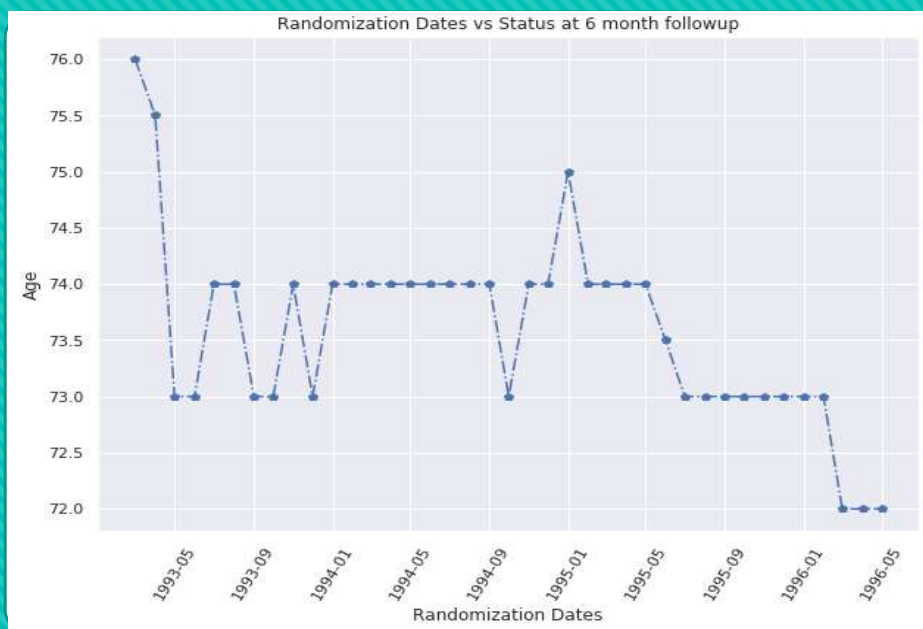
# Where do they end up?

# Blood pressure and Symptoms
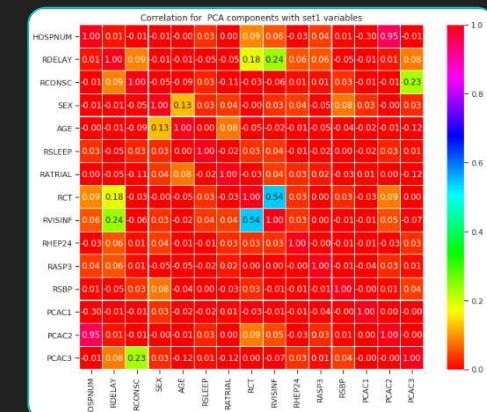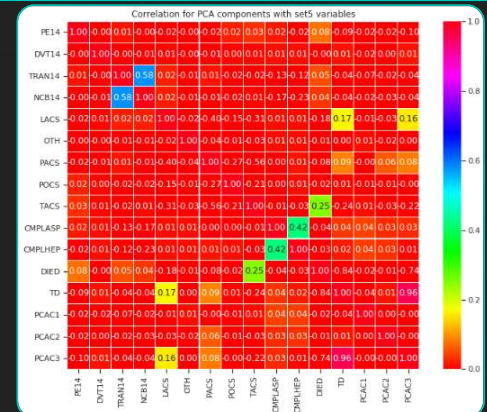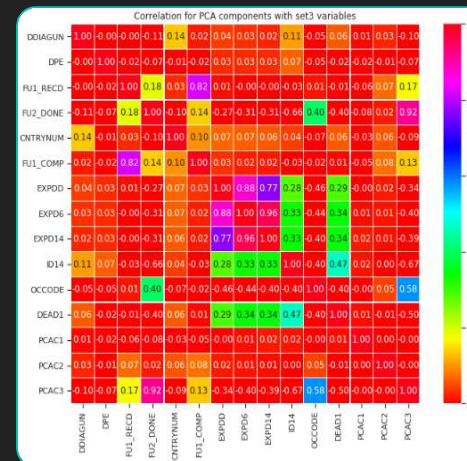
# Randomization data

# Age distribution and Atrial Fibrillation

# Unsupervised Modelling - PCA

Correlations of PCA components with other components in dataset containing patients 80 years or older.

# Supervised modelling

Let us try some of the supervised modelling techniques on patients 80 years and above and hence predict survival chances. Older people have been substantially under represented in stroke trials to date, so we hope the number of patients aged over 80 in this data set could also facilitate planning of trials in the 'older old'.

# Supervised Modelling Continued..

KNN Classifier Model

SVM Classifier Model

# Supervised Modelling Continued..

### SGD Classifier Model



### Decision Tree Classifier Model

# Supervised Modelling Continued..

### Random Forest Classifier Model



### ADA Boost Classifier Model

# Supervised Modelling Continued..

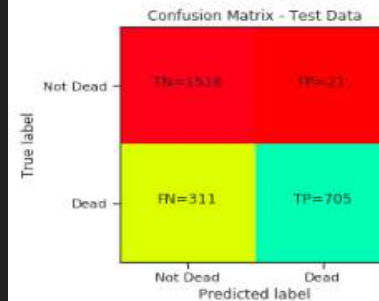### Gradient Boost Classifier Model



```
GradientBoostingClassifier
****Results****
Accuracy: 92.4794%
Cross validation scores: [0.95303327 0.962818   0.64187867 0.74168297 0.63850688]
Average accuracy of cross validation: 0.79 (+/- 0.29)
Log Loss: 0.2214369339579642
```

### Gaussian Naïve Bayes Classifier Model



```
GaussianNB
****Results****
Accuracy: 86.9957%
Cross validation scores: [0.92172211 0.95694716 0.94716243 0.89236791 0.84675835]
Average accuracy of cross validation: 0.91 (+/- 0.08)
Log Loss: 0.4124879295666476
```
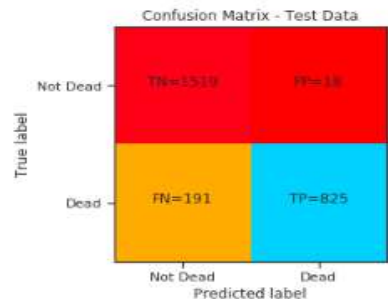
# Supervised Modelling Continued..

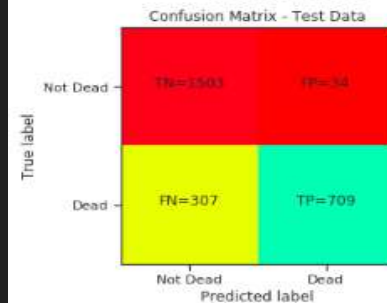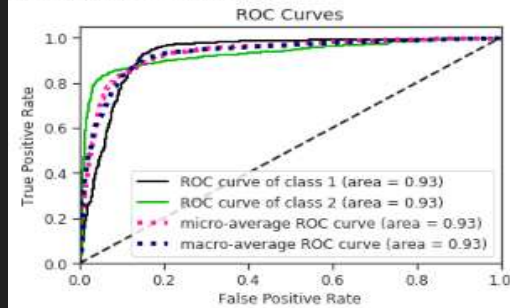### Linear Discriminant Analysis Model



### Quadratic Discriminant Analysis Model

# Best Model for patients 80 years and above

We are choosing Linear Discriminant Analysis as the best model in terms of prediction as well as log loss.

# Conclusion

Linear Discriminant Analysis was able to provide an average prediction of 85% survival chances for patients 80 years and above with an accuracy of 91% and a Loss of 25% with consistent Cross Validation Scores. Even though it appears that elderly have good survival chances, we need to keep in mind that our predictions were based on the dataset surveyed between 1993 to 1996.

Our dataset also doesn't capture some of the essential variables like height, weight, BMI, cholesterol, whether patients were smokers or alcoholics or diabetics and their measures, which, if provided, would help us make more accurate predictions.

Our predictions were for the 90s where medical facilities and surveying techniques were not as much advanced  as it is in the current age. Our next steps can be to survey the current stroke patients with all the essential variables and risk factors and use iteration and evaluation techniques and re-use the models which we used for this dataset and compare our findings and see if the differences were statistically significant or not using T-tests and P-Values.

# Next Steps

Prof. Anna Czlonkowska

Dr. hab Maciej Niewada

Prof. Peter Sandercock

# Credits

# Thank you!