# FINAL-TERM PROJECT REPORT

## Introduction To Data Science

Faculty
TOHEDUL ISLAM
Assistant Professor, Computer Science Dept.
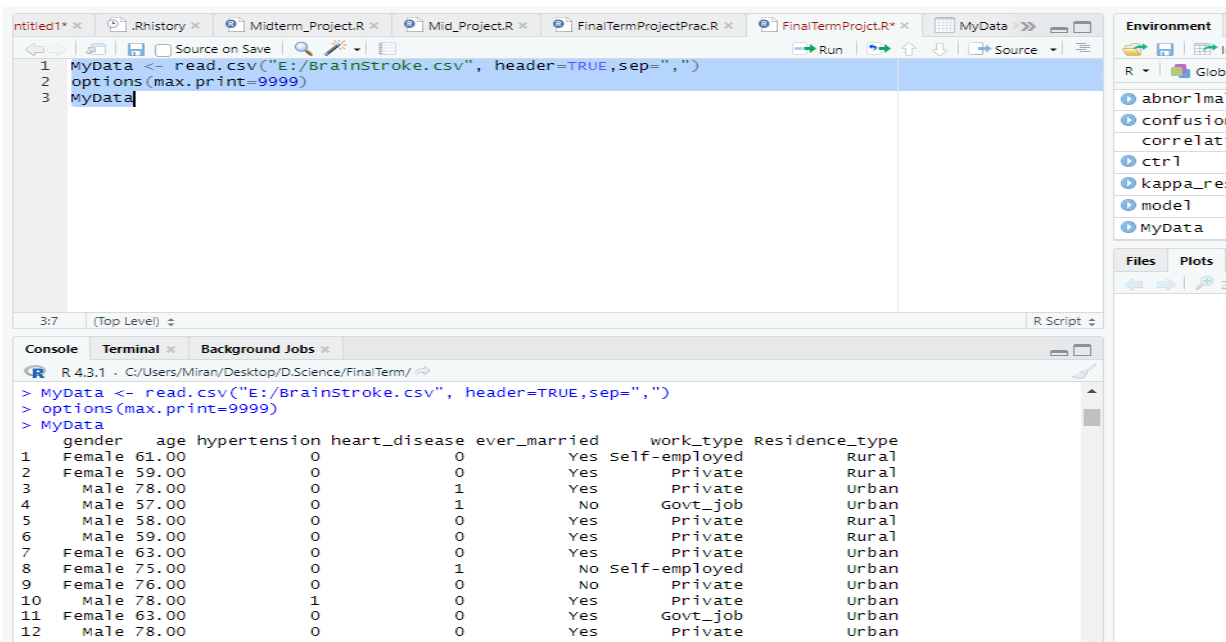
MD. ABDULLAH AL NASIR
19-41052-2

**About Dataset:** Here, the dataset named "Brain Stroke" has 202 instances and 11 attributes. Attributes are "Gender", "age", "hypertension", " heart_disease", " ever_married", " work_type", " Residence_type", " avg_glucose_level", " bmi", " smoking_status", " stroke".

Here, five attributes are categorical. Those are: "gender", "ever_married", "work_type", "Residence_type", "smoking_status". The attribute names: "age", "hypertension", " heart_disease", " avg_glucose_level", " bmi", " stroke" are numeric. The target attribute of this dataset is "stroke".

**Step- 01:** In the very first step, the dataset has been imported from the drive in a data frame named MyData using the read function. Using the options () function, we have imported all the data values with the help of max.print=9999.



**Step- 02:** The syntax of the function convertCategoricalToNumeric indicates the beginning of the code. This function expects MyData to be a data frame containing the dataset to be processed.

The function includes a list called categorical_cols that lists the column names that are required to be changed from categorical to numeric format. "Gender," "ever_married," "work_type," "Residence_type," and "smoking_status" are the columns listed here.

Each of the categorical columns listed in categorical_cols are iterated over in the code's loop. The actions listed below are carried out for each column:

a. **<u>Factor Conversion</u>**: The factor function converts the values of the column into factors. This step is essential because it assigns a numerical representation to each distinct value in the column.
b. **<u>Numeric Conversion</u>**: The as. numeric function is then used to convert the newly created factors into numeric values. This basically gives every individual factor value a numerical label.

**Step- 03:** The function colSums () is used to calculate the column sums of a matrix or data frame. The function is.na () is used to identify null values (NA) in a vector, matrix, or data frame. To calculate the number of missing values in each column, the function used here is colSums (is.na ()), keeping MyData as a parameter. Here we got that there are 0 null values in dataset.



**Step- 04:** Finding the outliers using Z-scores function. The scale () function is used to identify outliers by examining the standardized values of age, abs () retune the absolute value of a near. Outliers we found here are 1.32 3.00 6.00 7.00 0.48 5.00 1.88 1.08 1.80 3.00 1.32.

Now the outliers in the array named [age_outlier_indicates] are replaced with NA.



**Step- 05:** Another way to reduce outliers' values is to replace the most frequent or average value in the entire column. Using the table () function to calculate the most frequencies and which.max () to find the index of that value by using these functions together in the age attribute, we got the most frequent age, which is 78.

Now the data frame myFrequencyAgeData comparing with MyData and null values of age is replaced by 78.

**Step- 06:** Two additional columns are being added to MyData2: The values from the age column of the original MyData data frame are entered into the first column, ageRelation. The values from the stroke column of the initial MyData data frame are placed in the second column, strokeRelation.

To find the correlation coefficient between two variables, use the cor () function. In this instance, it will be applied to the MyData2 data frame's ageRelation and strokeRelation columns.

The variable correlation_coefficient contains the correlation calculation's outcome.



**Step- 07:** The R function cor () is used to find the correlation matrix between variables in a dataset.

MyData: It is believed that this is a dataset (data frame or matrix) with a number of different variables. The cor () function will determine how closely these variables are correlated.

"Pearson" is the correlation coefficient calculating method, according to this argument. The "Pearson" approach determines the Pearson correlation coefficient, which analyses the linear relationship between the variables.  Another R function, **round** (), rounds the values of a matrix or data frame to a predetermined number of digits.

digits = 2: This parameter indicates how many decimal places the correlation coefficients will be rounded to.

The correlation matrix will contain values between -1 and 1, where:

Values close to 1 indicate a strong positive correlation.

Values close to -1 indicate a strong negative correlation.

Values close to 0 indicate little or no linear correlation.

A strong positive correlation is found the output of this dataset.



**Step- 08:** The code corrplot (correlation_matrix, method = "color") generates a color-coded visualization of the correlation matrix stored in correlation_matrix. It uses the "corrplot" function, likely from the "corrplot" package, to create a plot where colors represent the strength and direction of correlations between variables.



**Step- 09:** install. packages("class") library(class) - The "class" package provides functions for k-nearest neighbors (KNN) classification and regression.

In this section we creating two data frame named training and test dataset. Here the code creates the train_data dataframe by extracting the "age" and "stroke" columns from the dataset, assuming there's one called MyData. The same columns are used to generate the test data dataframe as well.

The KNN algorithm's parameter k, which specifies how many neighbors to take into consideration when producing predictions, is set in this line. K in this case is set to 3.

The output starts with "0 1 1 1 1 1 1 1 ...", indicating the predicted class labels for the first few data points in the test dataset.

The output ends with "0 0 0 0 0 0 0", indicating the predicted class labels for the last few data points in the test dataset.



**Step- 10:** install. packages("caTools") library(caTools)- "caTools" package, which provides various functions for data manipulation and splitting datasets.

sample. split(...) - This function requires two arguments: the vector it want to split (in this case, the "stroke" column) and the SplitRatio parameter, which defines the percentage of data that should be given to the training set. In this case, the SplitRatio value of 0.8 indicates that around 80% of the data will be used for training and 20% for testing.

The train_data dataframe has a total of 161 rows, according to this output. Based on the outcome of the splitting procedure, which allocated about 80% of the data to the training set, these rows were chosen from the original dataset.

The test_data dataframe has 40 rows in total, according to this output. Based on the outcome of the splitting procedure, which allocated about 20% of the data to the test set, these rows were also chosen from the original dataset.

**Step − 11:** This code performs k-fold cross-validation for a k-nearest neighbors (KNN) classification model on the 'MyData' dataset. It starts by randomly assigning instances to 'num_folds' (which is 10 in this case) different folds. Then, for each fold, it separates the data into training and testing sets. The KNN algorithm is trained on the training set's 'age' and 'stroke' columns, and then predictions are made for the 'stroke' values in the test set. The accuracy of predictions for each fold is computed by comparing KNN predictions to the actual 'stroke' values in the test set. The accuracy scores for all folds are stored in the 'accuracy_scores' array. This process helps assess the KNN model's performance across different data subsets and provides insight into its overall effectiveness.



**Step-12:** By the help of the R function table (), one can count the instances of various combinations of values in two or more categorical variables.

Actual = test_data$stroke is specifying the values of the "Actual" column in the contingency table.

Predicted = knn_pred is specifying the values of the "Predicted" column in the contingency table.

```
80
81  install.packages("caret")
82  library(caret)
83
84
85  conf_matrix <- table(Actual = test_data$stroke, Predicted = knn_pred)
86  conf_matrix
87
6:12    (Top Level) ÷
```

```
nsole   Terminal ×   Background Jobs ×
  R 4.3.1 · C:/Users/Miran/Desktop/D.Science/FinalTerm/
66] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
vels: 0 1
library(caret)

conf_matrix <- table(Actual = test_data$stroke, Predicted = knn_pred)
conf_matrix
      Predicted
tual    0    1
   0  161    0
   1    6   34
```

**Step-13:** Precision value (1), indicating that when it predicts a positive instance (stroke), it's highly accurate. However, the recall value (0.85) suggests that the model is capturing most but not all of the actual positive instances.

This balance between precision and recall depends on the specific problem and the trade-offs you are willing to make between minimizing false positives and false negatives.

```
88
89  precision <- conf_matrix[2, 2] / sum(conf_matrix[, 2])
90  precision
91  recall <- conf_matrix[2, 2] / sum(conf_matrix[2, ])
92  recall
93
93:1    (Top Level) ÷
```

```
Console   Terminal ×   Background Jobs ×
 R  R 4.3.1 · C:/Users/Miran/Desktop/D.Science/FinalTerm/
1] "Precision: 1"
 print(paste("Recall:", recall))
1] "Recall: 0.85"
 precision <- conf_matrix[2, 2] / sum(conf_matrix[, 2])
 recall <- conf_matrix[2, 2] / sum(conf_matrix[2, ])
 precision <- conf_matrix[2, 2] / sum(conf_matrix[, 2])
 precision
1] 1
 recall
1] 0.85
```