

P8106_Midterm

jck2183_Chia-wen Kao

2021/3/26

```
library(tidyverse)
library(caret)
library(glmnet)
library(mlbench)
library(pROC) #generate ROC curve and calculate AUC
library(pdp) #partial dependent plot
library(vip) #variable importance plot: global impact on different predictor
library(AppliedPredictiveModeling) # for visualization purpose
library(corrplot)
library(RColorBrewer)
library(RANN)
library(visdat)
library(mgcv)
```

Introduction:

According to the World Health Organization (WHO) stroke is the 2nd leading cause of death globally, responsible for approximately 11% of total deaths. This dataset is used to predict whether a patient is likely to get stroke based on the input parameters like gender, age, various diseases, and smoking status. Each row in the data provides relevant information about the patient.

Data Source: <https://www.kaggle.com/fedesoriano/stroke-prediction-dataset>

All the features we had:

- id: unique identifier
- gender: “Male”, “Female” or “Other”
- age: age of the patient
- hypertension: 0 if the patient doesn’t have hypertension, 1 if the patient has hypertension
- heart_disease: 0 if the patient doesn’t have any heart diseases, 1 if the patient has a heart disease
- ever_married: “No” or “Yes”
- work_type: “children”, “Govt_joy”, “Never_worked”, “Private” or “Self-employed”
- Residence_type: “Rural” or “Urban”
- avg_glucose_level: average glucose level in blood
- bmi: body mass index
- smoking_status: “formerly smoked”, “never smoked”, “smokes” or “Unknown”*
- stroke: 1 if the patient had a stroke or 0 if not *Note: “Unknown” in smoking_status means that the information is unavailable for this patient

Import Data

```
stroke_df = read.csv("./data/healthcare-dataset-stroke-data.csv")
# head(stroke_df)

stroke_df$stroke = as.factor(stroke_df$stroke)
stroke_df$gender = factor(stroke_df$gender) %>% as.numeric()
stroke_df$ever_married = factor(stroke_df$ever_married) %>% as.numeric()
stroke_df$work_type = factor(stroke_df$work_type) %>% as.numeric()
stroke_df$Residence_type = factor(stroke_df$Residence_type) %>% as.numeric()
stroke_df$smoking_status = factor(stroke_df$smoking_status) %>% as.numeric()
stroke_df$heart_disease = factor(stroke_df$heart_disease) %>% as.numeric()
stroke_df$hypertension = as.numeric(factor(stroke_df$hypertension))
stroke_df$work_type = as.factor(stroke_df$work_type) %>% as.numeric()
stroke_df$bmi = as.numeric(stroke_df$bmi)
```

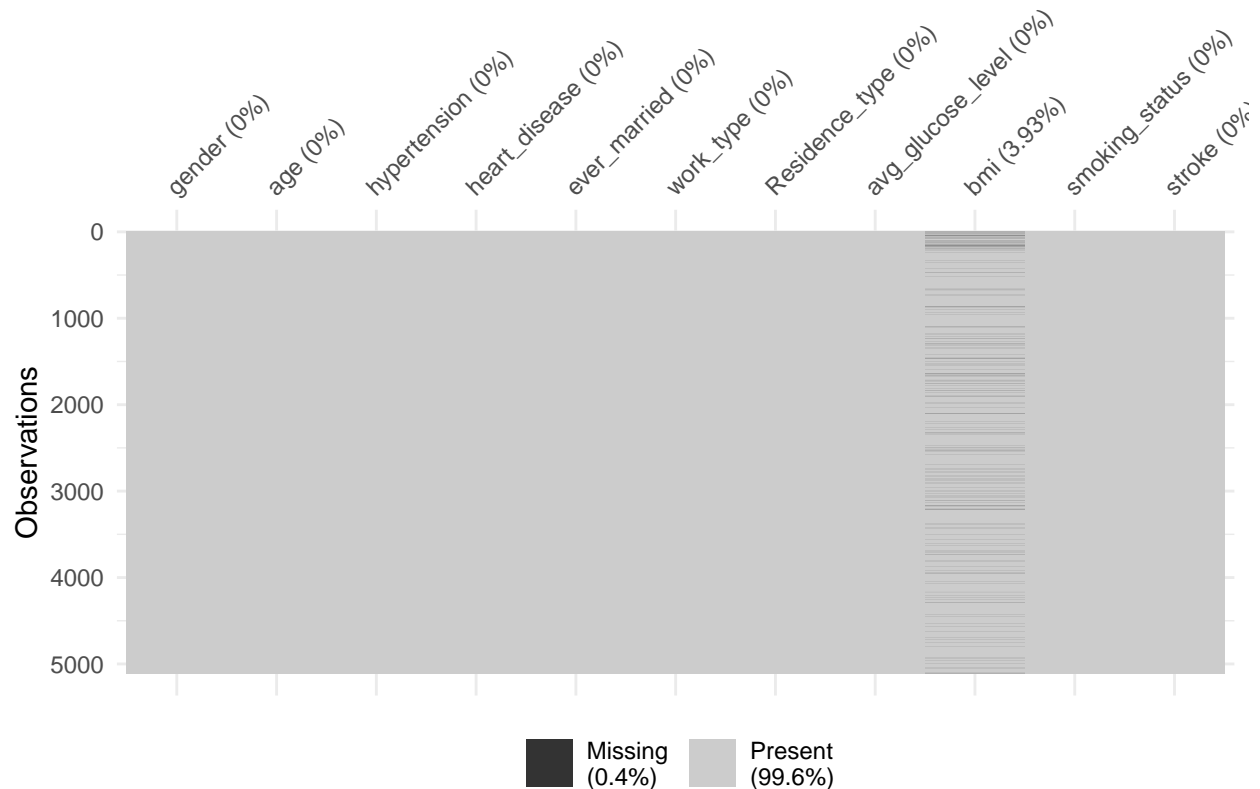
Warning: NAs introduced by coercion

```
stroke_df = stroke_df[, -1] %>%
  mutate(stroke = recode(stroke,
                        '0' = "No",
                        '1' = "Yes"),
         stroke = factor(stroke)) %>%
  filter(gender < 3)

summary(stroke_df)
```

```
##      gender      age      hypertension      heart_disease
## Min.   :1.000   Min.   : 0.08   Min.   :1.000   Min.   :1.000
## 1st Qu.:1.000   1st Qu.:25.00   1st Qu.:1.000   1st Qu.:1.000
## Median :1.000   Median :45.00   Median :1.000   Median :1.000
## Mean   :1.414   Mean   :43.23   Mean   :1.097   Mean   :1.054
## 3rd Qu.:2.000   3rd Qu.:61.00   3rd Qu.:1.000   3rd Qu.:1.000
## Max.   :2.000   Max.   :82.00   Max.   :2.000   Max.   :2.000
##
##      ever_married      work_type      Residence_type      avg_glucose_level
## Min.   :1.000   Min.   :1.000   Min.   :1.000   Min.   : 55.12
## 1st Qu.:1.000   1st Qu.:2.000   1st Qu.:1.000   1st Qu.: 77.24
## Median :2.000   Median :4.000   Median :2.000   Median : 91.88
## Mean   :1.656   Mean   :3.495   Mean   :1.508   Mean   :106.14
## 3rd Qu.:2.000   3rd Qu.:4.000   3rd Qu.:2.000   3rd Qu.:114.09
## Max.   :2.000   Max.   :5.000   Max.   :2.000   Max.   :271.74
##
##      bmi      smoking_status      stroke
## Min.   :10.30   Min.   :1.000   No :4860
## 1st Qu.:23.50   1st Qu.:2.000   Yes: 249
## Median :28.10   Median :2.000
## Mean   :28.89   Mean   :2.586
## 3rd Qu.:33.10   3rd Qu.:4.000
## Max.   :97.60   Max.   :4.000
## NA's      :201
```

```
vis_miss(stroke_df)
```



The imported dataset has 5110 observations in total. Excluding the id, we only gave ten features and one binary outcome variable-stroke (0:no stroke, 1:stroke). We found that the stroke outcome distribution is imbalanced with 4861 observations have no stroke while 249 observations have a stroke.

We find out there are 201 observations with missing values in BMI. Among these missing values, 40 observations have a stroke while 161 observations without stroke. We will then apply preprocess imputation in the caret train function to address the imputation problem. We also have 1544 unknown in smoke status, will treat those who answered unknown as a variable so no need to impute them.

Our main task is to find out the appropriate models that have a better performance on prediction by comparing several models' performance.

First, we have to convert character variables into factors to add them into our model and proceed with the analysis. Plus, we will also examine if there is any correlation among features. Meanwhile, we also found there is an observation who identified their gender as "Other". We decide to omit this single subject so that we can proceed with our analysis.

Next, the characteristics of features will help us determine which model would be proper. As the outcome is binary, and the features are mixtures of continuous and categorical variables. We also have to decide how to partition the train and test data, which cross-validation method to use. Evaluation metrics should be used and set up a reasonable tuning grid corresponding to the tuning parameter.

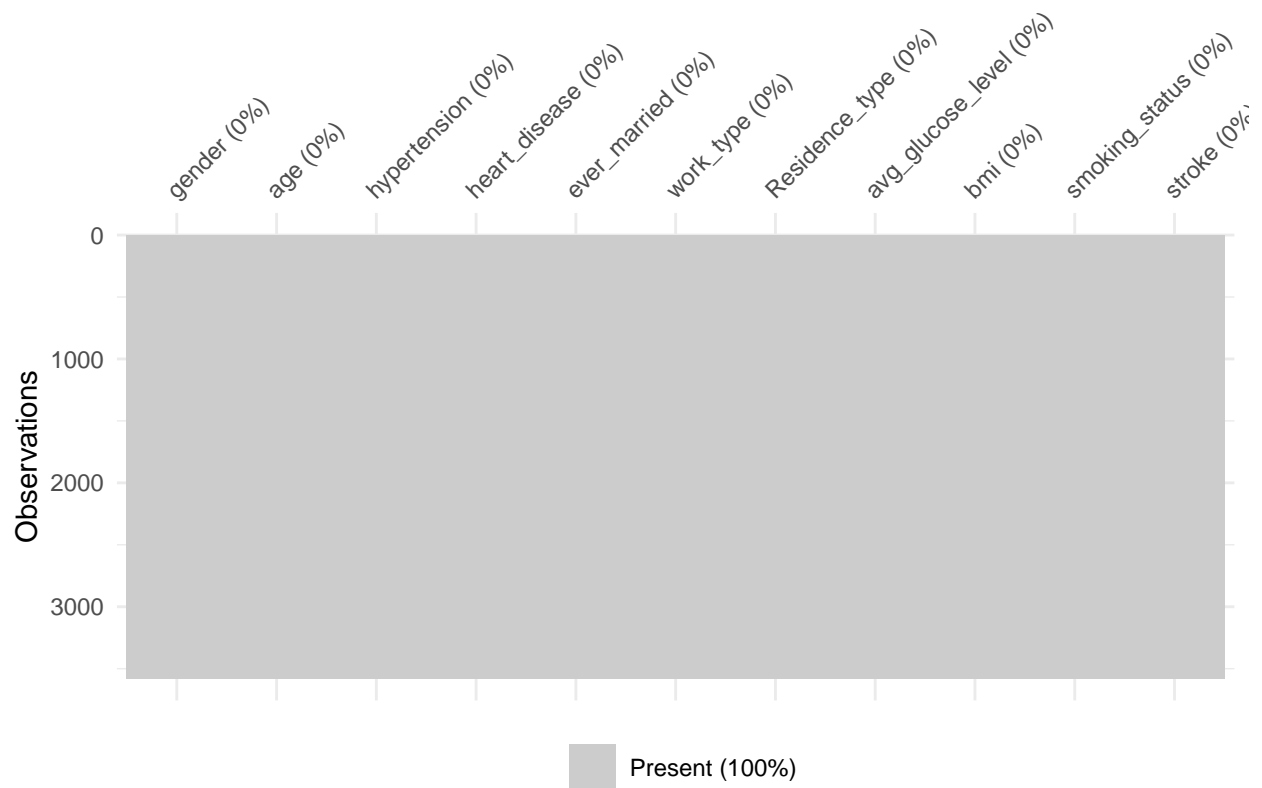
Exploratory Data Analysis

Partition the dataset, I will use 70% as training data and 30% as test data.

```
set.seed(123)
trRow = createDataPartition(y = stroke_df$stroke, p = 0.7, list = F)
train.data = stroke_df[trRow, ]
test.data = stroke_df[-trRow, ]
```

Try imputation with `preProcess()`

```
knnImp = preProcess(train.data, method = "knnImpute", k = 3)
train.data = predict(knnImp, train.data)
vis_miss(train.data)
```



```
test.data = predict(knnImp, test.data)
vis_miss(test.data)
```

