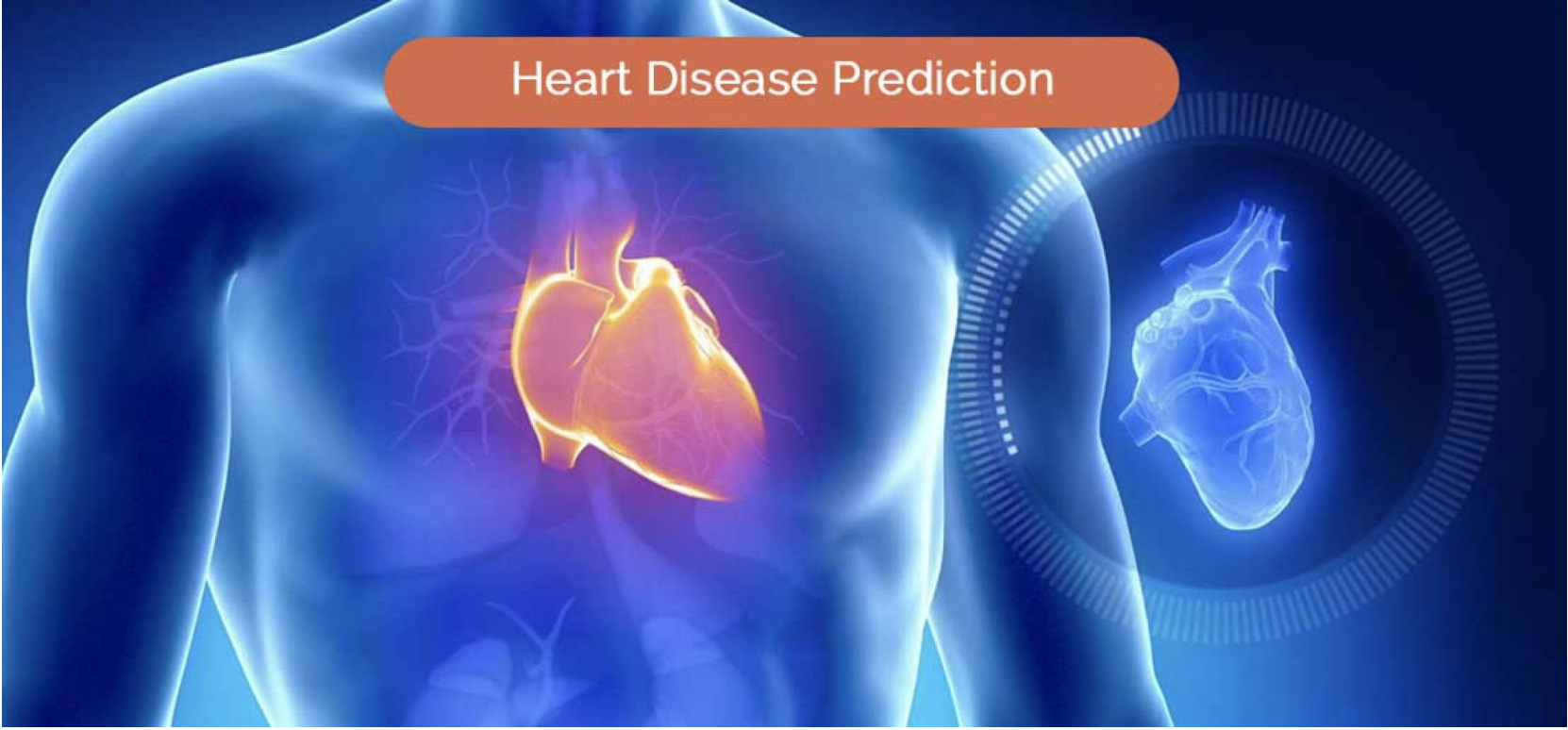


# HEART DISEASE PREDICTION

Ruitong Liu, Ying Liu, Chaoying Luo, Bo Yang, Yiyang Sheng

## Background

World Health Organization has estimated 12 million deaths occur worldwide, every year due to Heart diseases. Half the deaths in the United States and other developed countries are due to cardiovascular diseases. The early prognosis of cardiovascular diseases can aid in making decisions on lifestyle changes in high-risk patients and in turn reduce the complications. This project intends to pinpoint the most relevant/risk factors of heart disease as well as predict the overall risk using logistic regression, decision tree and random forest.

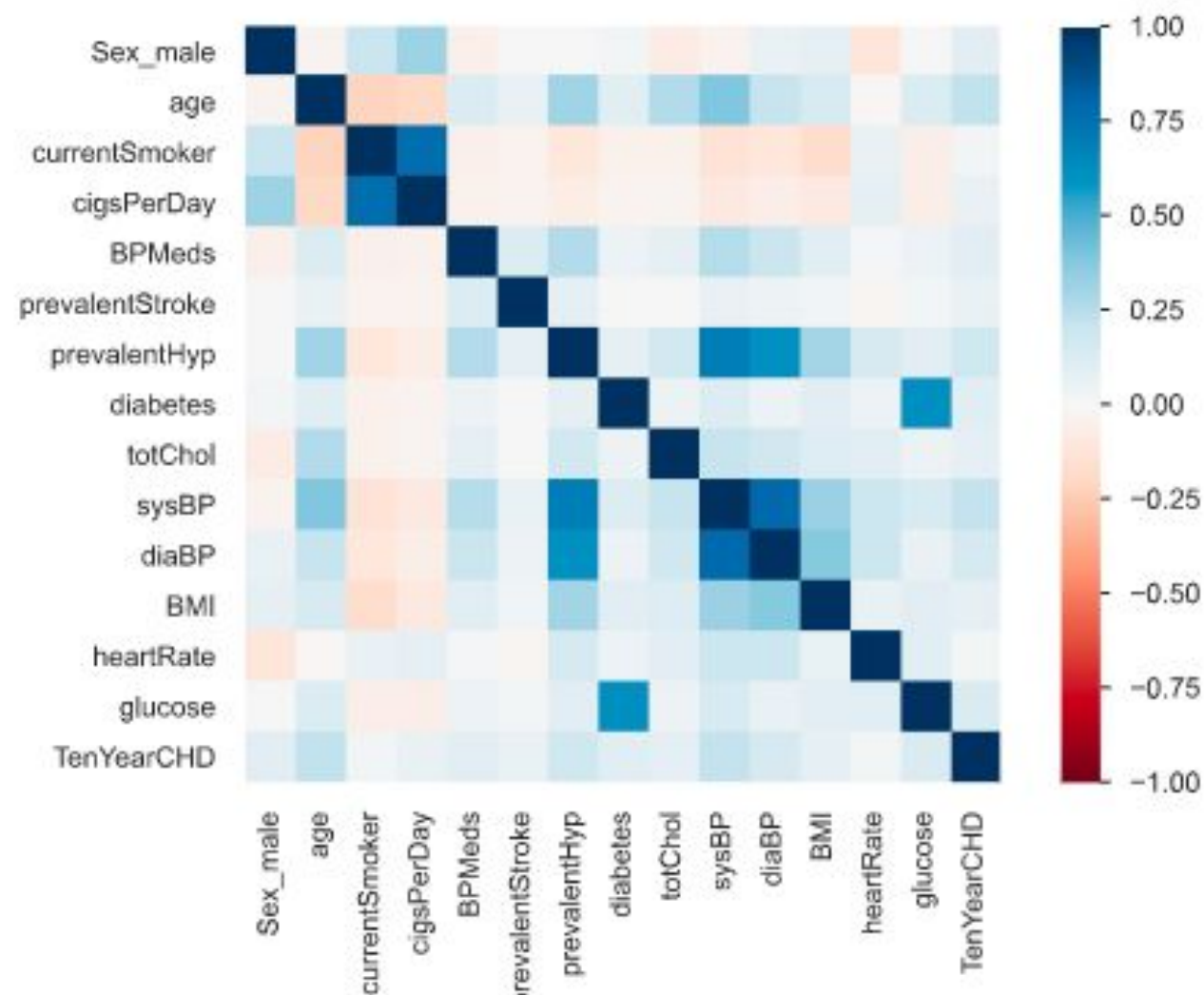
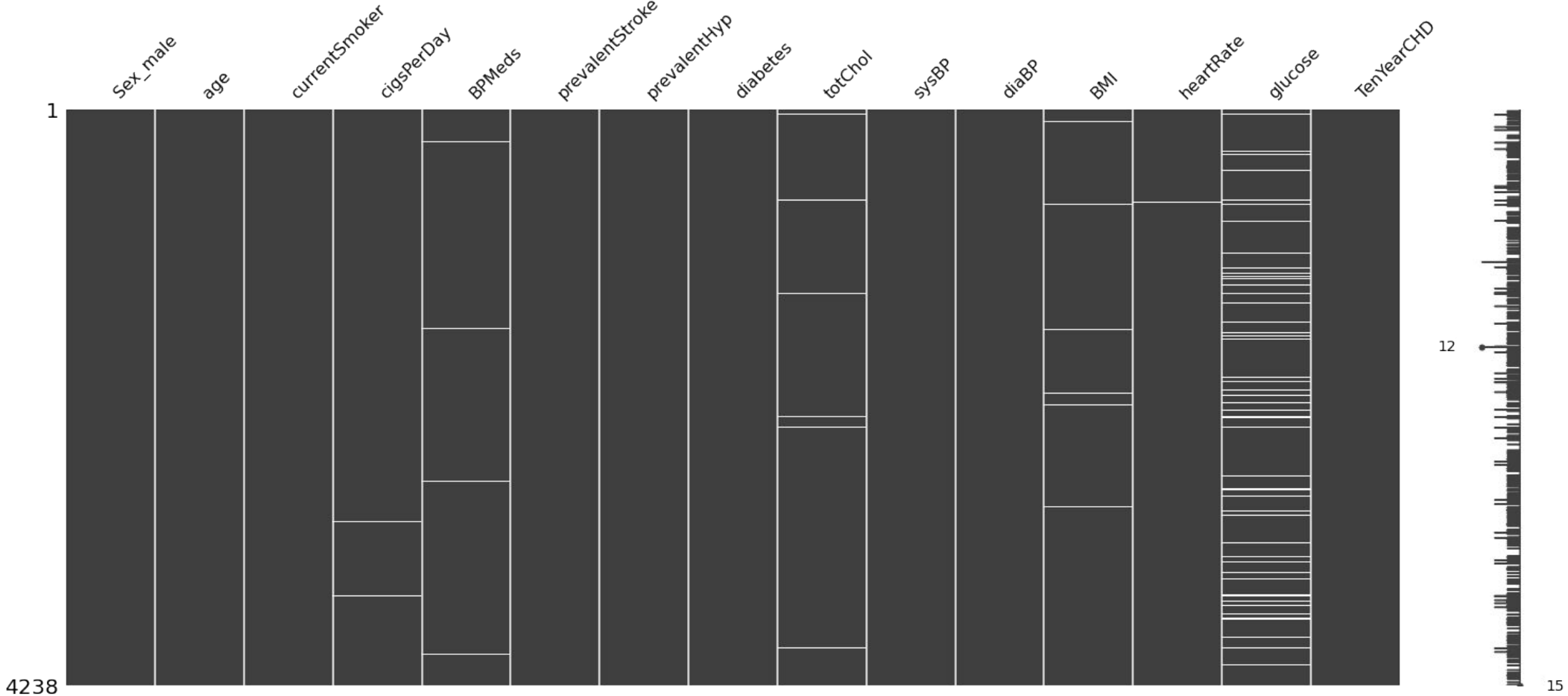


## Dataset

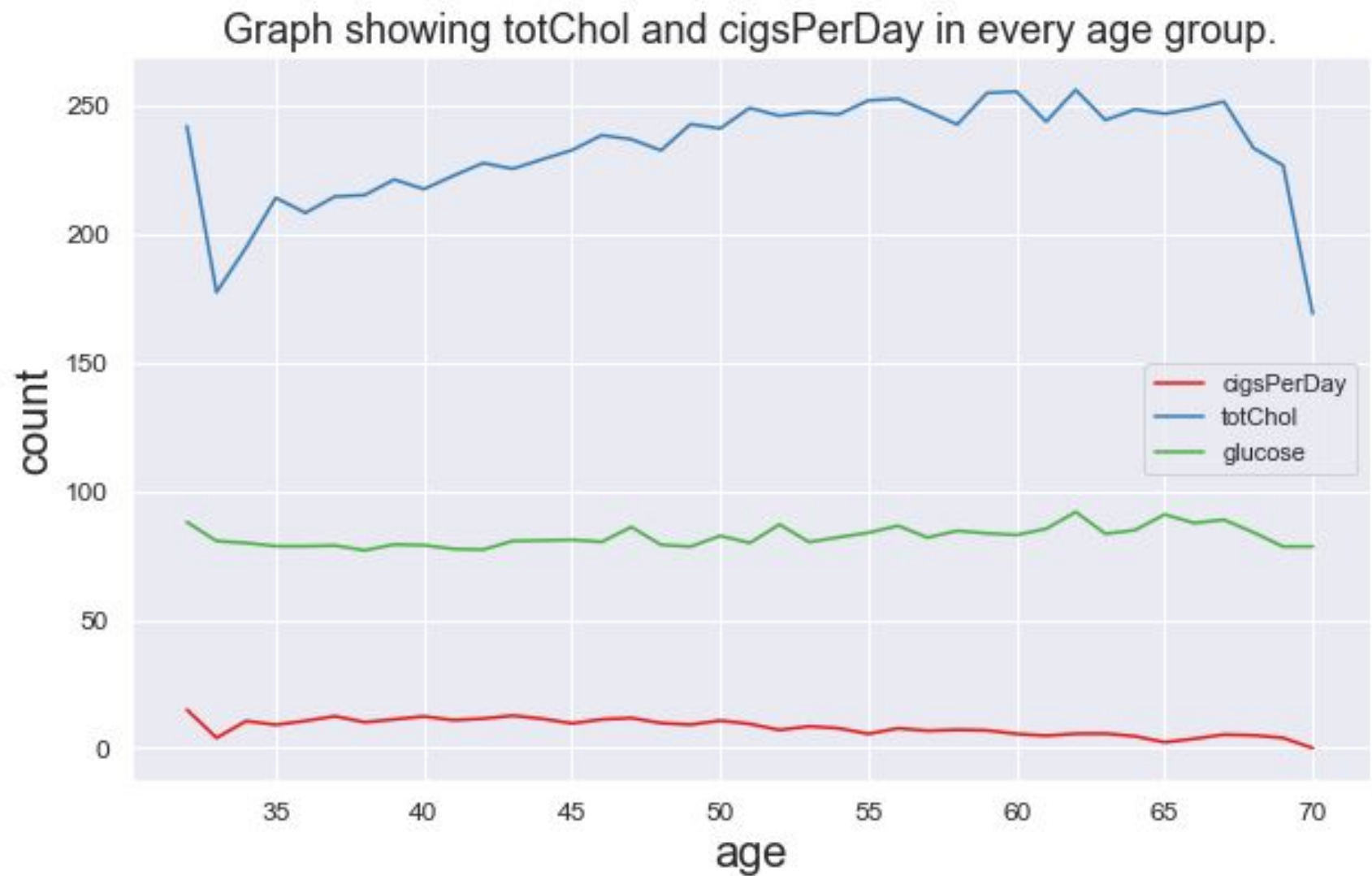
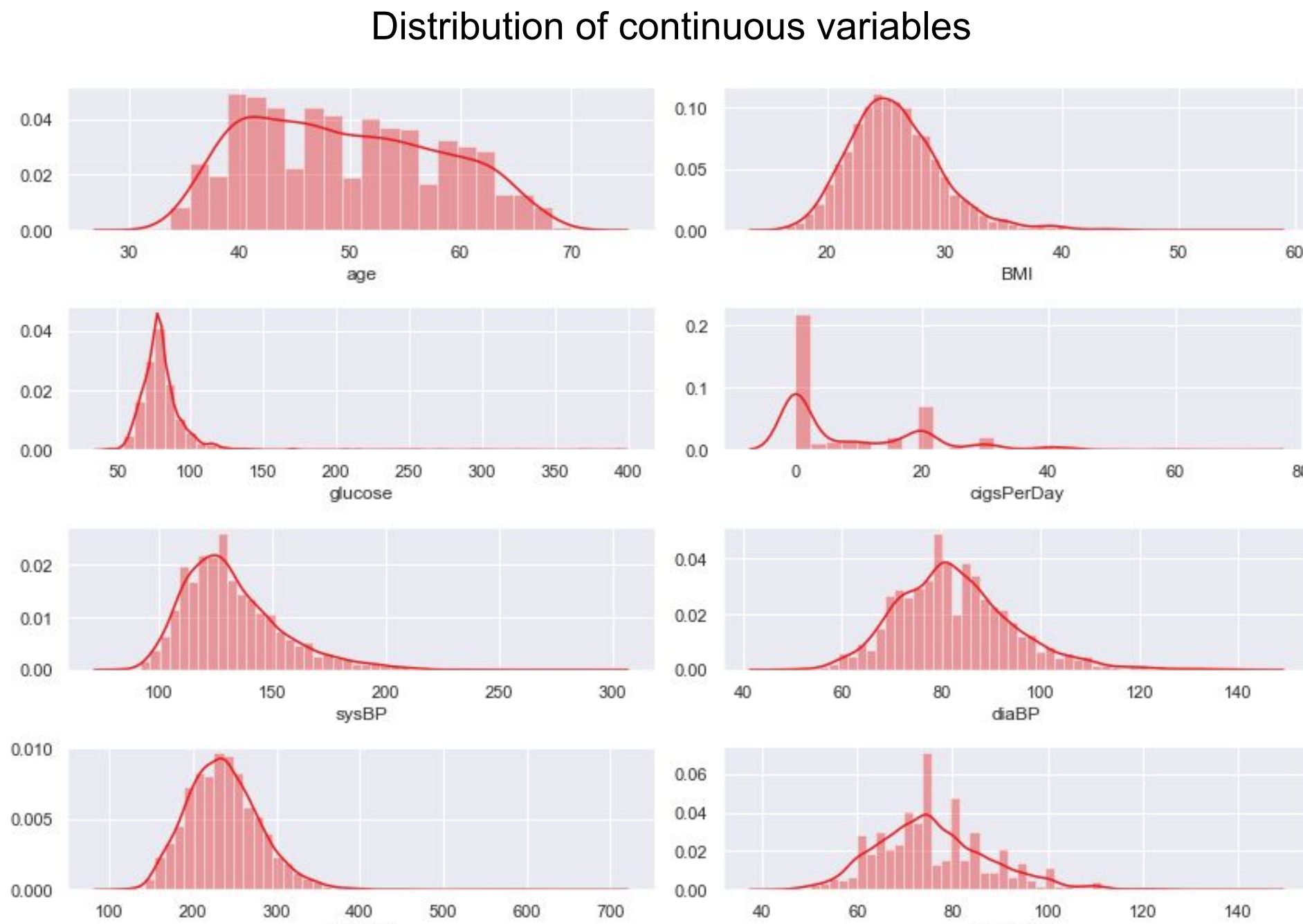
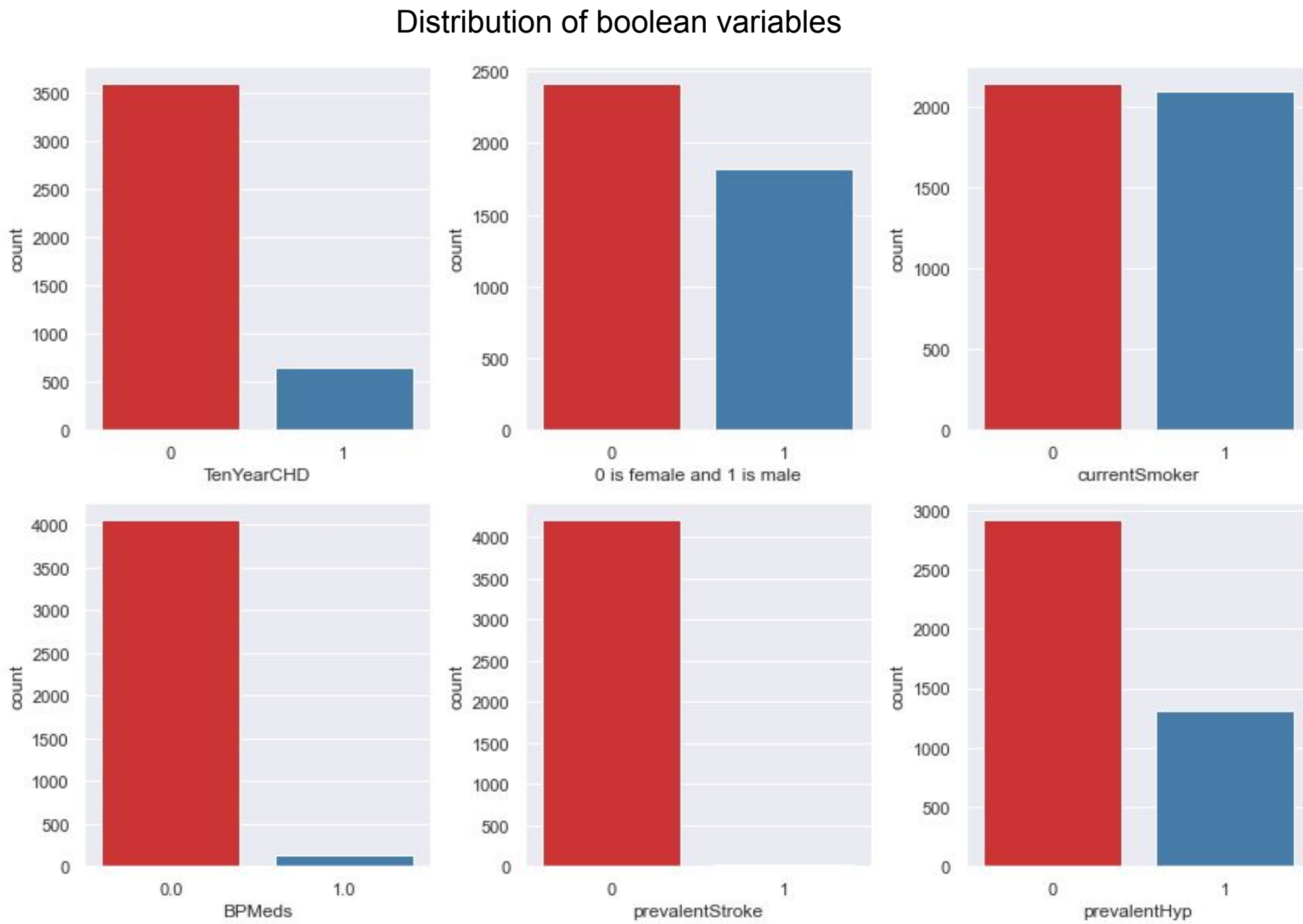
The dataset is publicly available on the Kaggle website. The classification goal is to predict whether the patient has 10-year risk of future coronary heart disease (CHD). The dataset provides the patients' information. It includes over 4,000 records and 15 attributes. Each attribute is a potential risk factor. There are both demographic, behavioral and medical risk factors.

- Sex
- Age
- Current Smoker: whether or not the patient is a current smoker
- Cigs Per Day: the number of cigarettes that the person smoked on average in one day.
- BP Meds: whether or not the patient was on blood pressure medication
- Prevalent Stroke: whether or not the patient had previously had a stroke
- Prevalent Hyp: whether or not the patient was hypertensive
- Diabetes: whether or not the patient had diabetes
- Tot Chol: total cholesterol level
- Sys BP: systolic blood pressure
- Dia BP: diastolic blood pressure
- BMI: Body Mass Index
- Heart Rate: heart rate
- Glucose: glucose level
- 10 year risk of coronary heart disease CHD

## Data Cleaning



## EDA



## Logistic Regression

```
Call:
glm(formula = TenYearCHD ~ ., family = "binomial", data = training_data)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.3368  -0.5844  -0.4280  -0.2956   2.8490

Coefficients:
(Intercept)      -8.420864    0.768718 -10.954 < 2e-16 ***
Sex_male         -0.529831    0.119547   4.432 9.34e-06 ***
age              -0.059878    0.007341   8.157 3.44e-16 ***
currentSmoker1   -0.105741    0.175227   0.603 0.546289
cigsPerDay       -0.025283    0.006907   3.269 0.001078 **
BPmeds1         -0.388775    0.263768   1.474 0.140489
prevalentStroke1 1.004468    0.559026   1.806 0.070480 .
prevalentHyp1    -0.138832    0.154691   0.892 0.372227
diabetes1        -0.165525    0.368816  -0.449 0.653575
totChol          -0.001152    0.001221   0.944 0.345178
sysBP            -0.016231    0.004214   3.852 0.000117 ***
diaBP            -0.000538    0.007092  -0.075 0.940425
BMI              -0.018621    0.013977  -0.768 0.447311
heartRate        -0.001806    0.004615   0.218 0.827415
glucose          -0.018481    0.002663   3.936 8.30e-05 ***

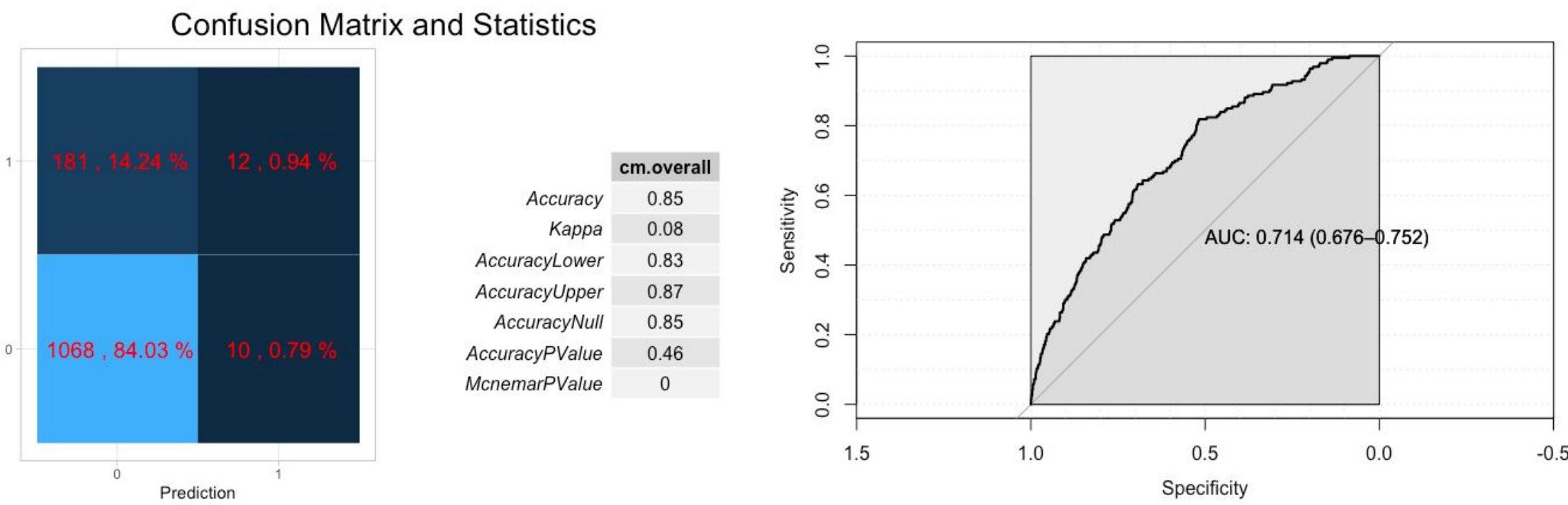
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

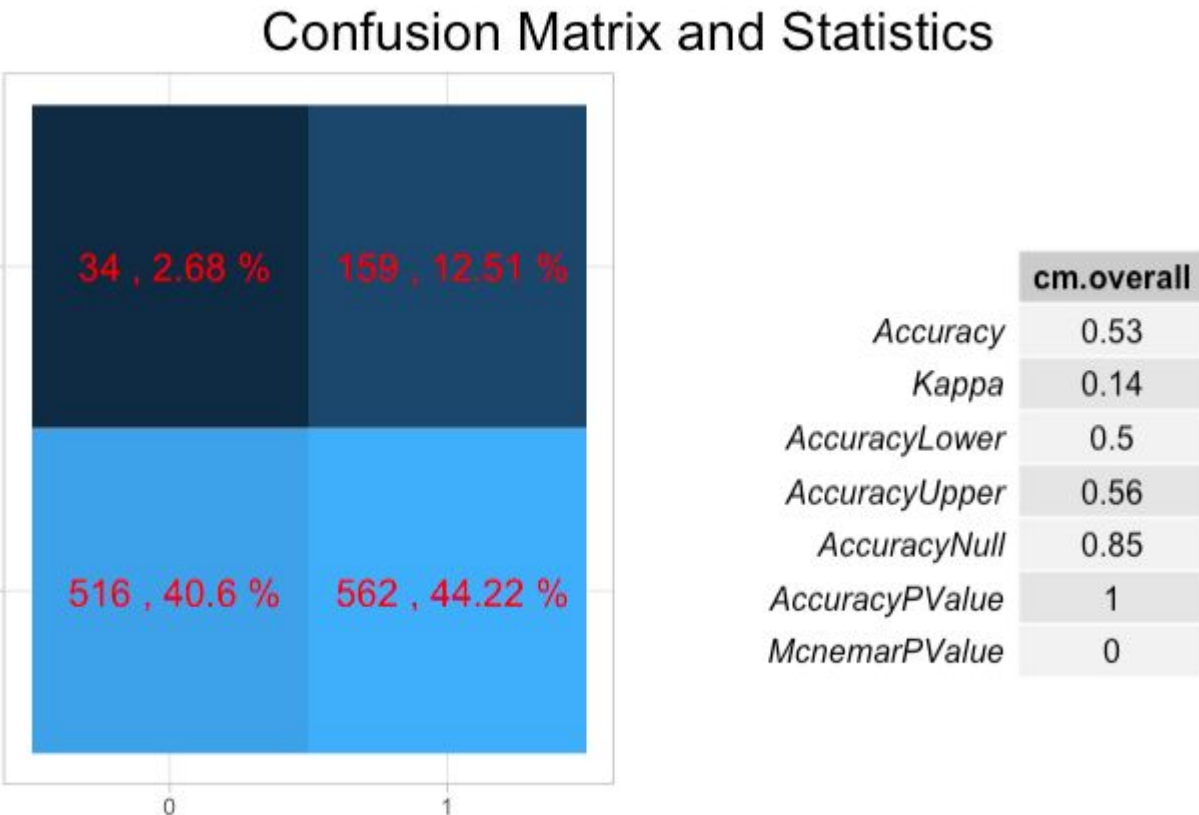
Null deviance: 2528.9 on 2966 degrees of freedom
Residual deviance: 2236.9 on 2952 degrees of freedom
AIC: 2250.1

Number of Fisher Scoring iterations: 5
```

Significant variables include Sex\_male, age, cigsPerDay, prevalentStroke, sysBP, and glucose. We can run logitic regression again with these significant variables.



Since the model is predicting Heart disease too many type II errors is not advisable. A False Negative (ignoring the probability of disease when there actually is one) is more dangerous than a False Positive in this case. Hence in order to increase the sensitivity, threshold can be lowered. 0.1 threshold is adopted.

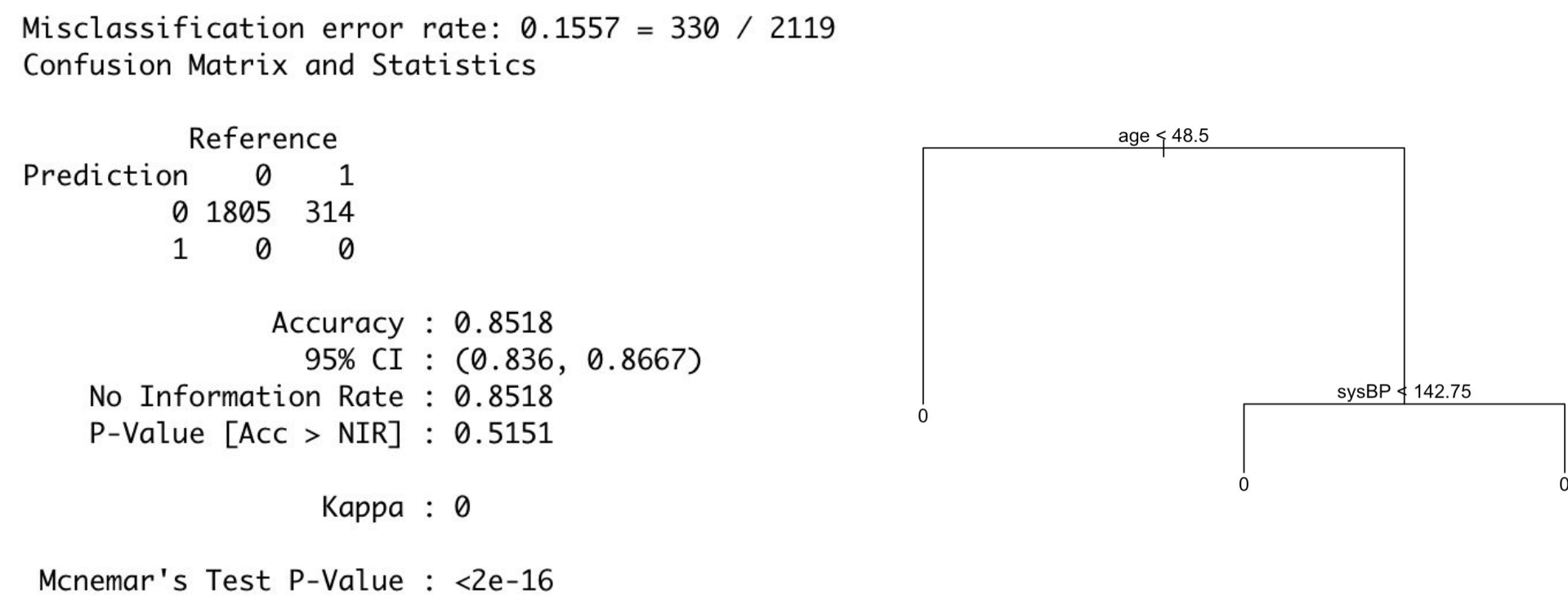
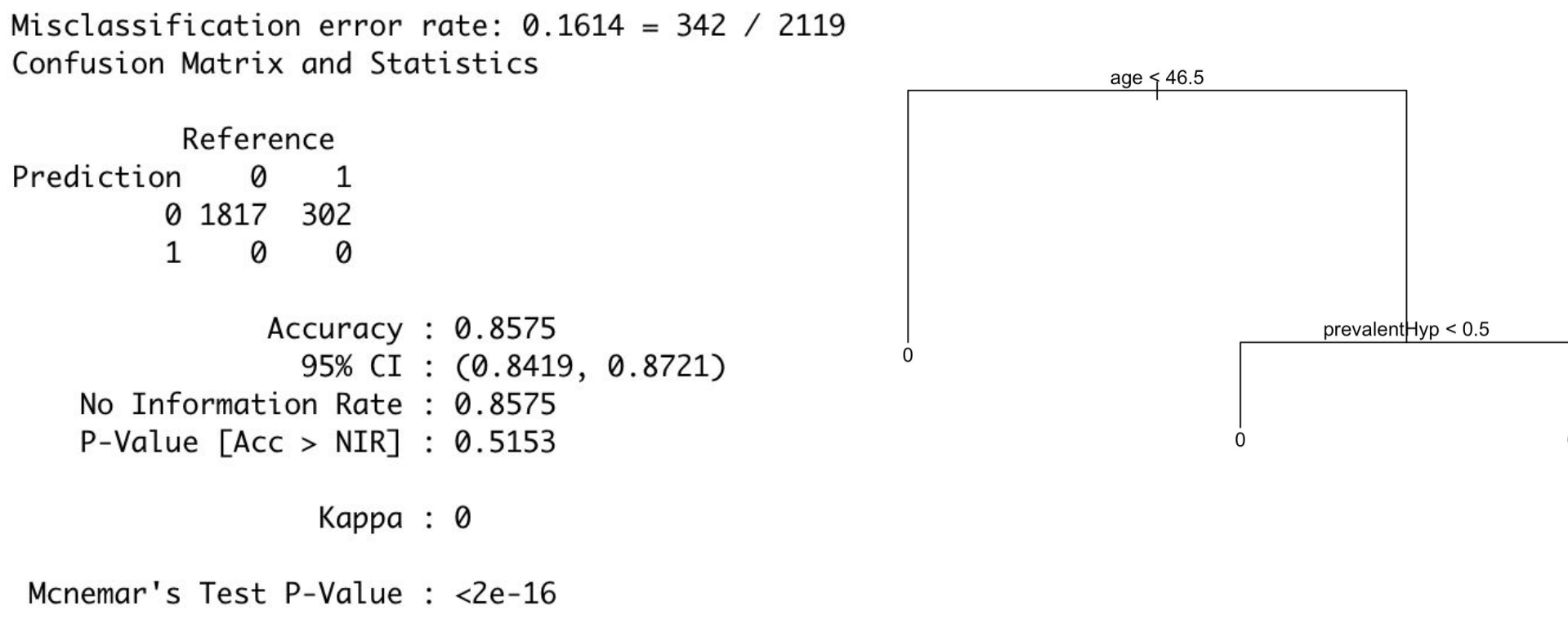


Now as we have a good accuracy for predicting people with chances of heart disease in next 10 years, we can treat them and guide them well in advance.

## Decision Tree

Building a decision tree with all the variables, we have an accuracy around 85.75% with a 16.14% misclassification error rate. And the variables included are age and prevalentHyp. If we use the significant variables from logistic regression, the predict accuracy is 85.18% with a 15.57% misclassification error rate. And the variables included are age and sysBP. although the first tree gave us a slightly better prediction accuracy, the second tree has a lower error rate and all the variables are significant variables. Therefore, the second tree has a better performance.

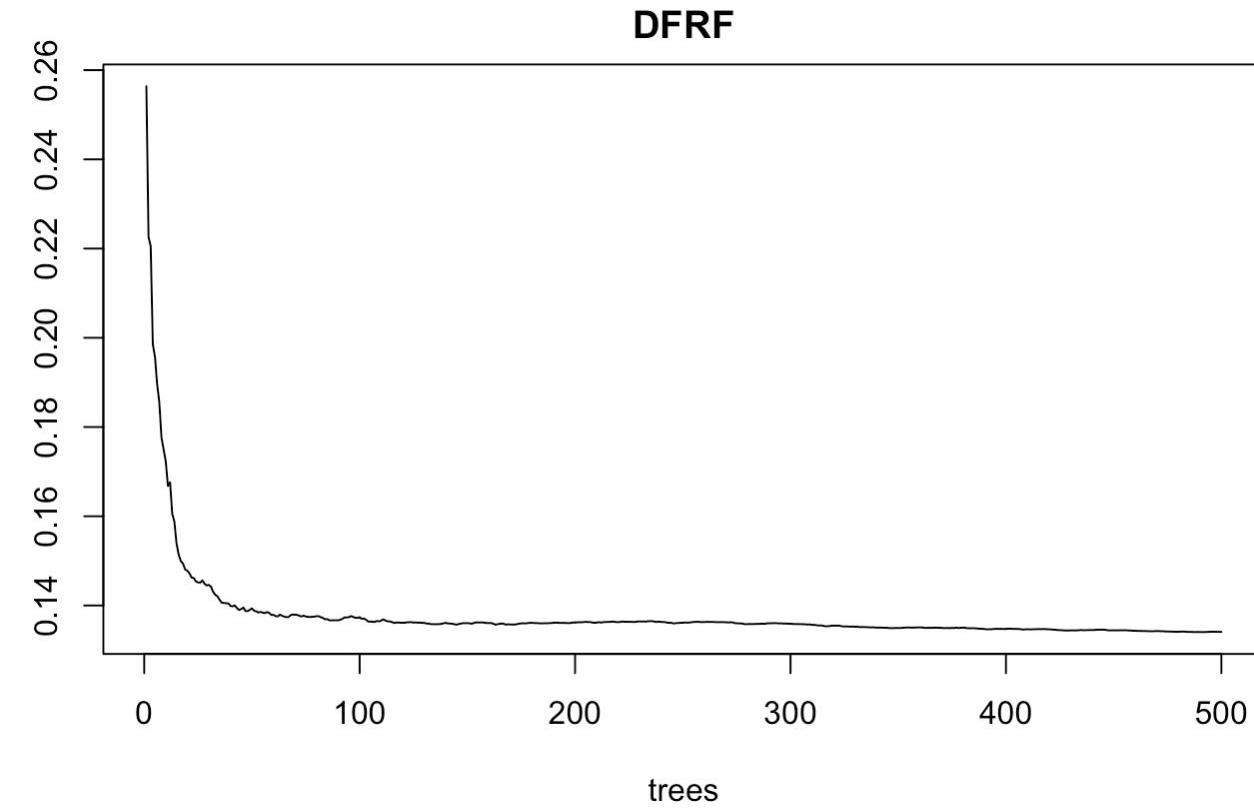
## Decision Tree



## Random Forest

age	11.6070253	16.7266999
currentSmoker	2.6635266	1.2461663
cigsPerDay	2.2489689	6.4644841
prevalentHyp	7.6938274	2.9433389
diabetes	2.5322369	0.4492975
totChol	1.0916327	17.0078529
BMI	1.8202488	20.5723299
heartRate	-1.0991501	13.4242186
glucose	-0.6624642	14.9095989

The importance decreases when "MeanDecreaseAccuracy" in the table above decreases. So set "MeanDecreaseAccuracy" in decreasing order, then it shows that "age", "prevalentHyp" and "currentSmoker" are the top 3 significant predictors.



The error is minimized at roughly 100 trees.

## Conclusion

- In terms of accuracy, the best model is the decision tree. In terms of variable abundance, the best model is the logistic regression.
- The patient that was hypertensive or had diabetes shows increasing odds of having heart disease. Increase in age, number of cigarettes smoked per day and systolic Blood Pressure also show increasing odds of having heart disease.
- Men seem to be more susceptible to heart disease than women.
- Meanwhile, people under the age of 48.5 have a very low probability of suffering from heart disease.