



## Multi Class Prediction of Obesity Risk



## Project Overview

- Developing a model which will accurately predict obesity risk among individuals based on various lifestyle factors

# Dataset

- **id:** A unique identifier for each individual in the dataset.
- **Gender:** The individual's gender, indicating whether they are male or female.
- **Age:** The age of the individual, representing their age in years.
- **Height:** The height of the individual, typically measured in meters.
- **Weight:** The weight of the individual, typically measured in kilograms.
- **family\_history\_with\_overweight:** Indicates whether there is a family history of overweight for the individual (yes/no).
- **FAVC:** Stands for "Frequency of consuming high caloric food," representing how often the individual consumes high-calorie foods (yes/no).
- **FCVC:** Stands for "Frequency of consuming vegetables," representing how often the individual consumes vegetables.
- **NCP:** Stands for "Number of main meals," indicating the number of main meals the individual consumes daily.
- **CAEC:** Stands for "Consumption of food between meals," representing the frequency of consuming food between meals.
- **SMOKE:** Indicates whether the individual smokes or not (yes/no).
- **CH2O:** Represents the amount of water consumption for the individual.
- **SCC:** Stands for "Calories consumption monitoring," indicating whether the individual monitors their calorie consumption (yes/no).
- **FAF:** Stands for "Physical activity frequency," representing the frequency of the individual's physical activities.
- **TUE:** Stands for "Time using technology devices," indicating the amount of time the individual spends using technology devices.
- **CALC:** Stands for "Consumption of alcohol," representing the frequency of alcohol consumption.
- **MTRANS:** Stands for "Mode of transportation," indicating the mode of transportation the individual uses.
- **NObesidad:** The target variable, representing the obesity risk category of the individual. It has multiple classes such as 'Overweight\_Level\_II', 'Normal\_Weight', 'Insufficient\_Weight', 'Obesity\_Type\_III', 'Obesity\_Type\_II', 'Overweight\_Level\_I', and 'Obesity\_Type\_I'.

# Pre-Processing



Overviewing the Dataset



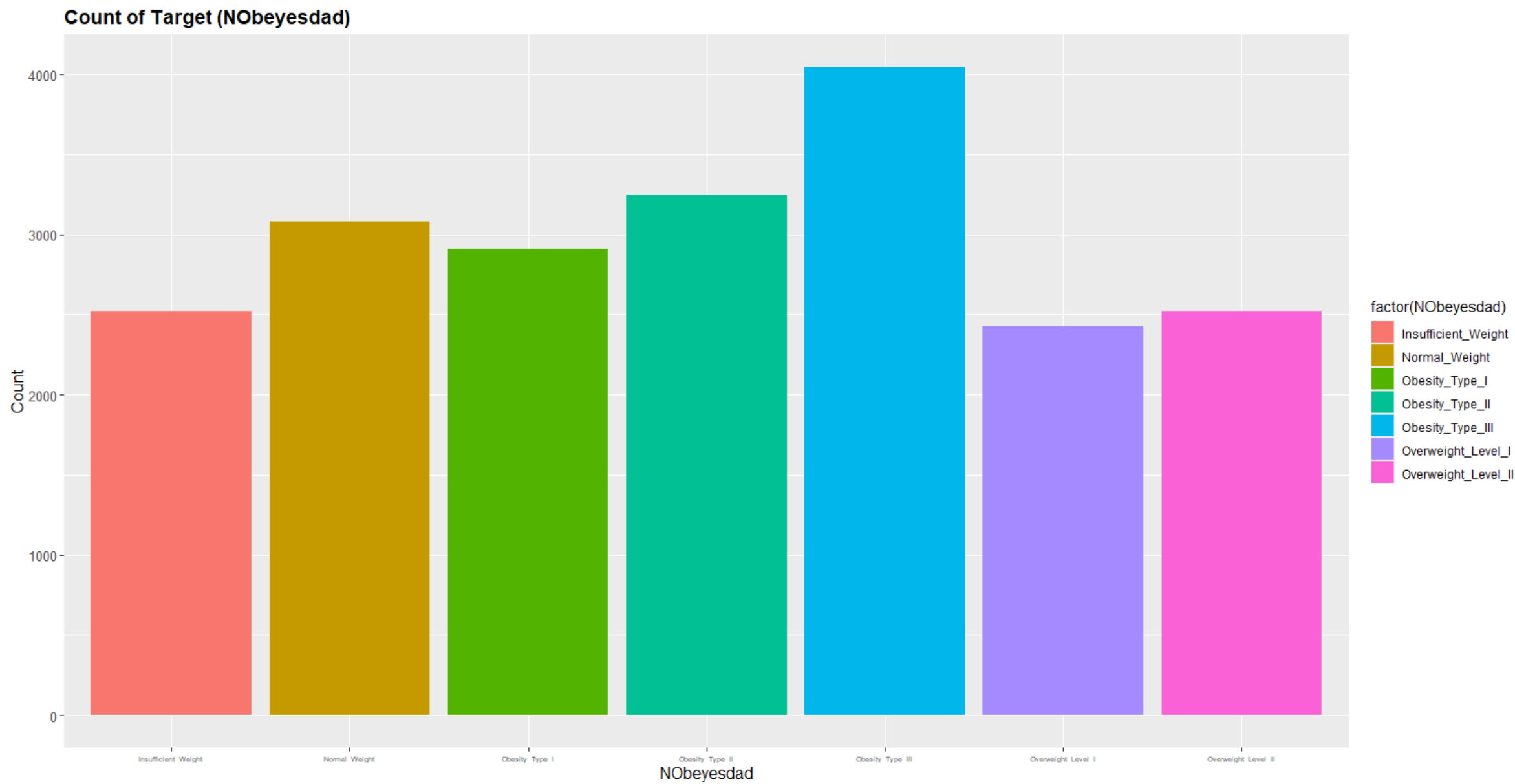
Missing Values  
and duplicates



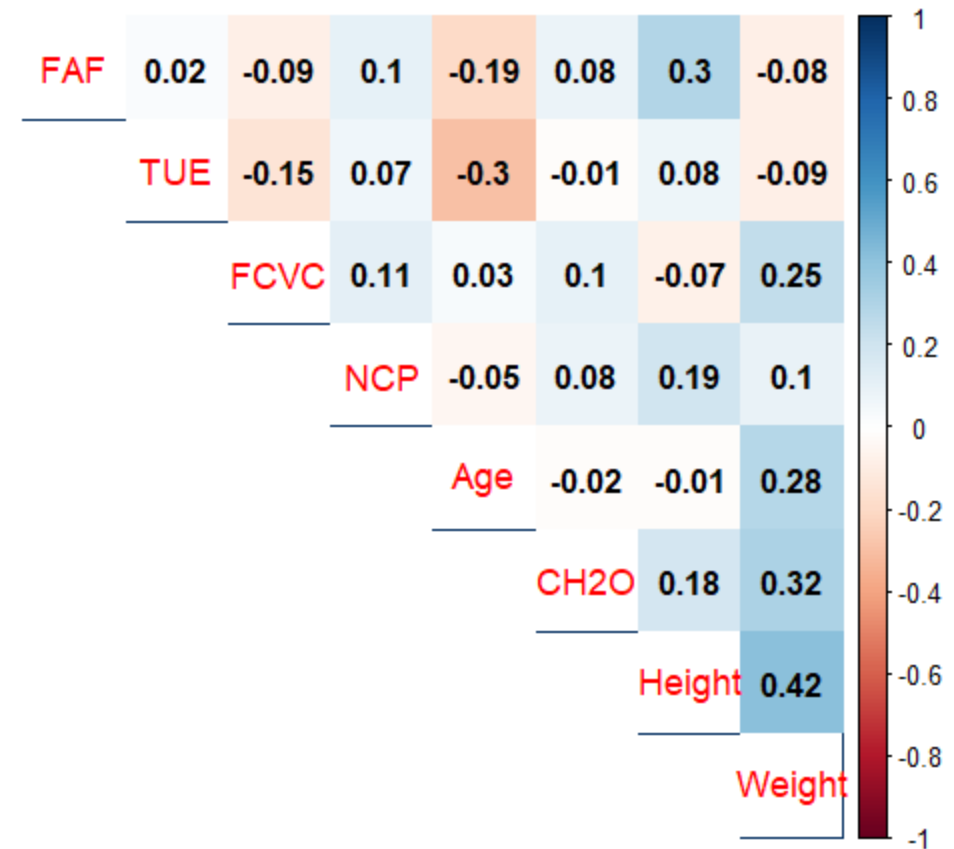
Encoding all variables into  
numerical variables



Feature Engineering



# Correlation between numerical variables



# Missing Values and Duplicates

```
##           Gender           Age
##           0              0
##           Height          Weight
##           0              0
## family_history_with_overweight  FAVC
##           0              0
##           FCVC            NCP
##           0              0
##           CAEC            SMOKE
##           0              0
##           CH2O            SCC
##           0              0
##           FAF             TUE
##           0              0
##           CALC            MTRANS
##           0              0
##           NObeyesdad
##           0
```

```
# duplicates across train/ test
sum(
  duplicated(rbind(df_train[, -17], df_test))
)
```

```
## [1] 0
```

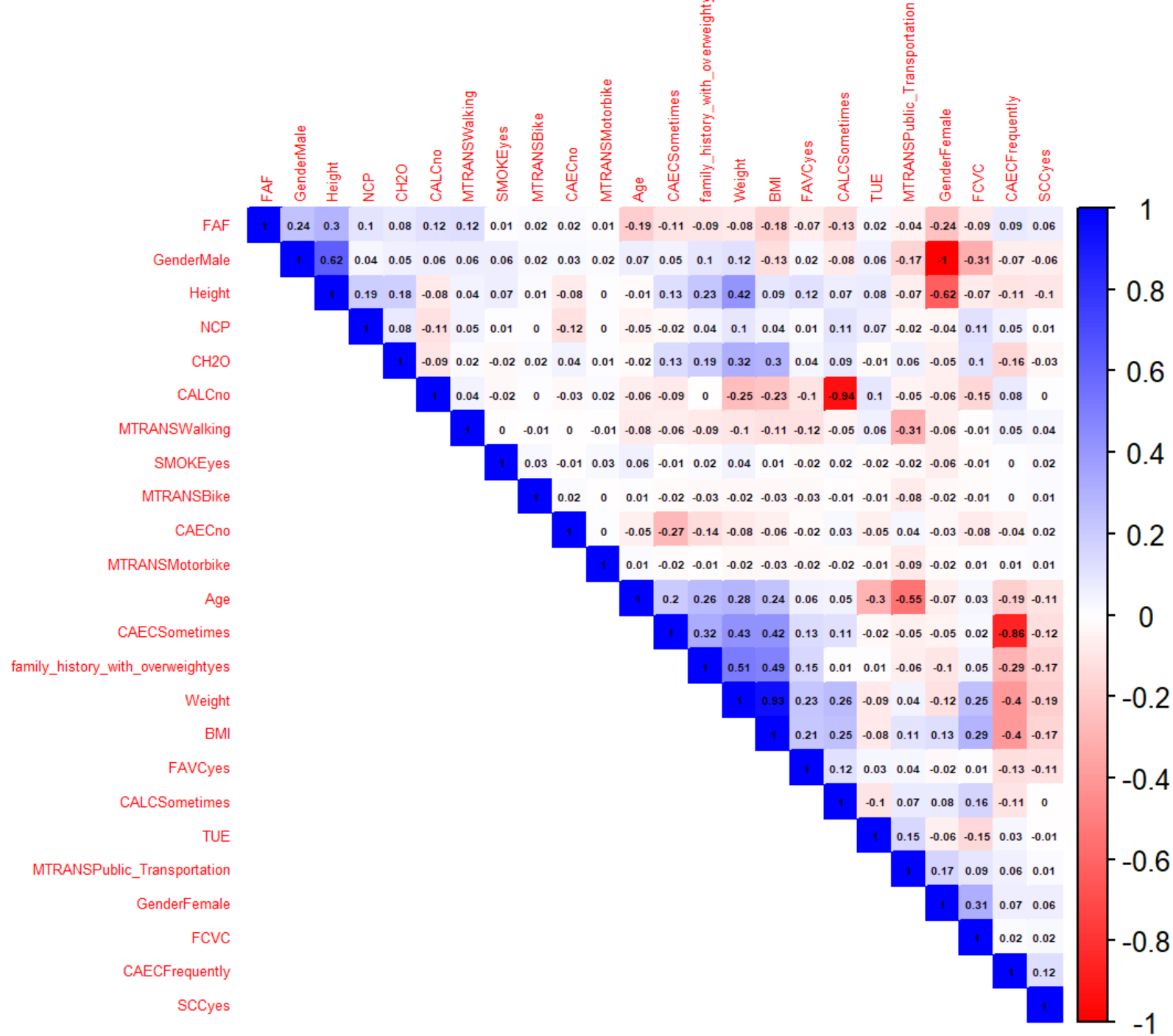
```
# null values by column
colSums(is.na(df_train))
```

# Adding BMI column adjusted by age and gender

- `calculate_BMI <- function(weight, height, age, gender_male) {`
- `bmi <- weight / (height ^ 2)`
- 
- `# Adjust BMI for age and gender`
- `if (gender_male == 1) {`
- `if (age < 18) {`
- `bmi <- bmi * 1.1`
- `} else if (age >= 18 & age <= 24) {`
- `bmi <- bmi * 1.05`
- `}`
- `} else {`
- `if (age < 18) {`
- `bmi <- bmi * 1.15`
- `} else if (age >= 18 & age <= 24) {`
- `bmi <- bmi * 1.08`
- `}`
- `}`

```
$ GenderFemale : num [1:20758] 0 1 1 1 0 0 0 0
$ GenderMale : num [1:20758] 1 0 0 0 1 1 1 1
$ Age : num [1:20758] 24.4 18 18 21 31
$ Height : num [1:20758] 1.7 1.56 1.71 1.
$ weight : num [1:20758] 81.7 57 50.2 131
$ family_history_with_overweightyes : num [1:20758] 1 1 1 1 1 1 1 1
$ FAVCyes : num [1:20758] 1 1 1 1 1 1 1 1
$ FCVC : num [1:20758] 2 2 1.88 3 2.68
$ NCP : num [1:20758] 2.98 3 1.41 3 1.
$ CAECFrequently : num [1:20758] 0 1 0 0 0 0 0 0
$ CAECno : num [1:20758] 0 0 0 0 0 0 0 0
$ CAECSometimes : num [1:20758] 1 0 1 1 1 1 1 1
$ SMOKEyes : num [1:20758] 0 0 0 0 0 0 0 0
$ CH2O : num [1:20758] 2.76 2 1.91 1.67
$ SCCyes : num [1:20758] 0 0 0 0 0 0 0 0
$ FAF : num [1:20758] 0 1 0.866 1.468
$ TUE : num [1:20758] 0.976 1 1.674 0.
$ CALCno : num [1:20758] 0 1 1 0 0 0 0 0
$ CALCSometimes : num [1:20758] 1 0 0 1 1 1 1 1
$ MTRANSBike : num [1:20758] 0 0 0 0 0 0 0 0
$ MTRANSMotorbike : num [1:20758] 0 0 0 0 0 0 0 0
$ MTRANSPublic_Transportation : num [1:20758] 1 0 1 1 1 1 0 0
$ MTRANSWalking : num [1:20758] 0 0 0 0 0 0 0 0
$ NObeyesdad : Factor w/ 7 levels "Insufficie
$ BMI : num [1:20758] 28.3 25.3 18.5 4
```

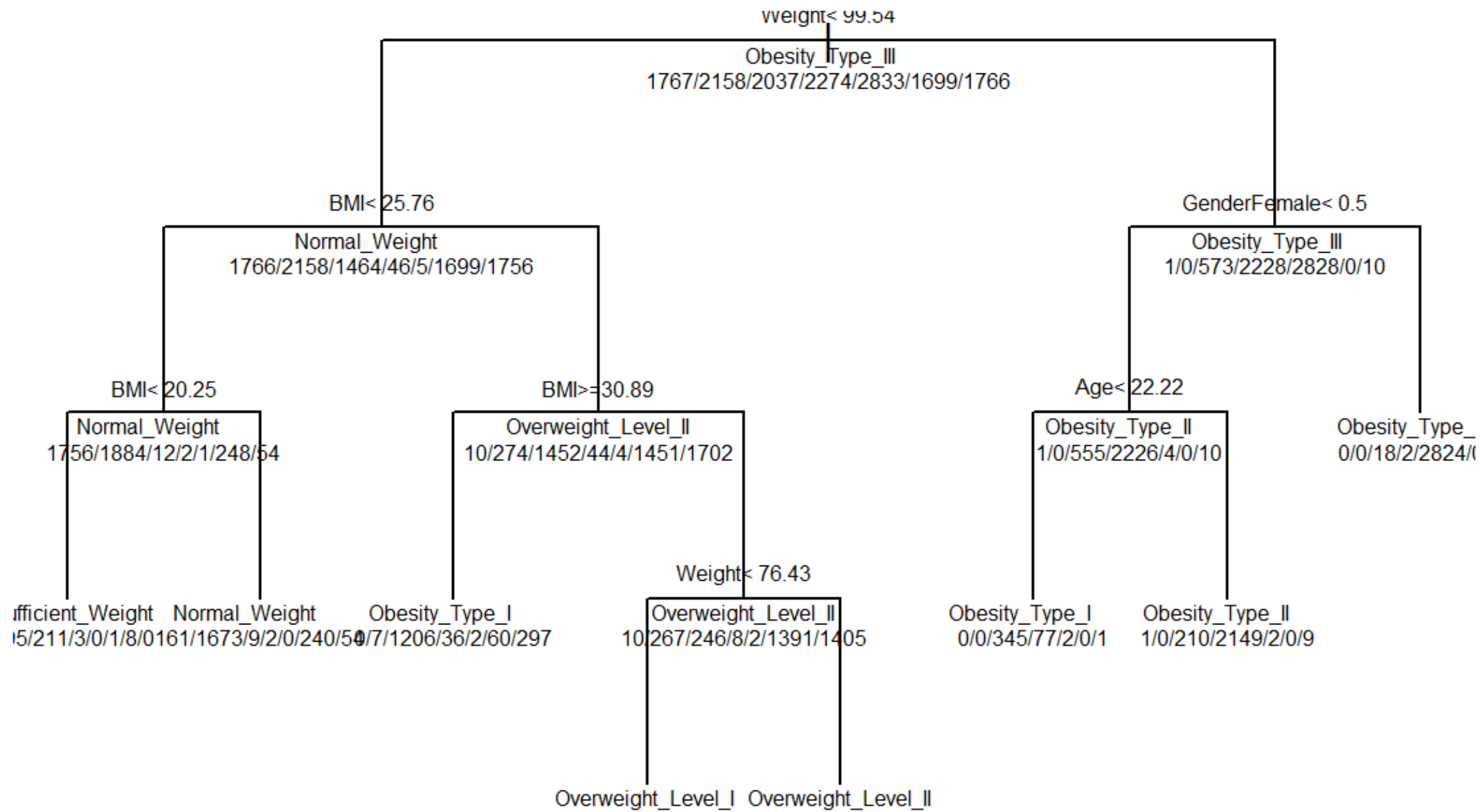




# Models used to compute output

- Decision tree
- Logistic regression
- Random forest

## Decision Tree for Obesity Prediction



# Decision tree

- Confusion Matrix and Statistics

- 

- Reference

Prediction	Insufficient_Weight	Normal_Weight	Obesity_Type_I	Obesity_Type_II	Obesity_Type_III	Overweight_Level_I	Overweight_Level_II
• Insufficient_Weight	673	75	0	0	0	6	0
• Normal_Weight	77	753	3	0	0	87	16
• Obesity_Type_I	0	1	643	52	0	21	117
• Obesity_Type_II	0	0	92	920	0	0	7
• Obesity_Type_III	0	0	10	1	1211	0	1
• Overweight_Level_I	5	92	29	0	1	563	144
• Overweight_Level_II	1	3	96	1	1	51	471

## Overall Statistics

Accuracy : 0.8409

95% CI : (0.8316, 0.8499)

No Information Rate : 0.1949

P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.8134

Mcnemar's Test P-Value : NA

# Logistic Regression

Prediction	Insufficient_Weight	Normal_Weight	Obesity_Type_I	Obesity_Type_II	Obesity_Type_III	Overweight_Level_I	Overweight_Level_II
<b>Insufficient_Weight</b>	693	97	0	0	0	8	0
<b>Normal_Weight</b>	59	757	2	0	0	74	16
<b>Obesity_Type_I</b>	0	0	714	37	1	13	79
<b>Obesity_Type_II</b>	0	0	47	928	1	0	6
<b>Obesity_Type_III</b>	0	0	6	0	1210	0	1
<b>Overweight_Level_I</b>	2	62	27	0	1	513	103
<b>Overweight_Level_II</b>	2	8	77	9	0	120	551

Overall Statistics Accuracy : 0.8621

95% CI : (0.8533, 0.8706)

No Information Rate : 0.1949

P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.8383

Mcnemar's Test P-Value : NA

# Random Forest

Prediction	Insufficient_Weight	Normal_Weight	Obesity_Type_I	Obesity_Type_II	Obesity_Type_III	Overweight_Level_I	Overweight_Level_II
<b>Insufficient_Weight</b>	697	41	0	0	0	6	0
<b>Normal_Weight</b>	53	825	4	0	1	85	17
<b>Obesity_Type_I</b>	0	1	774	14	1	18	48
<b>Obesity_Type_II</b>	0	0	30	951	0	0	7
<b>Obesity_Type_III</b>	0	0	4	1	1211	0	0
<b>Overweight_Level_I</b>	5	44	11	0	0	530	50
<b>Overweight_Level_II</b>	1	13	50	8	0	89	634

## Overall Statistics

Accuracy : 0.9026

95% CI : (0.8965, 0.9113)

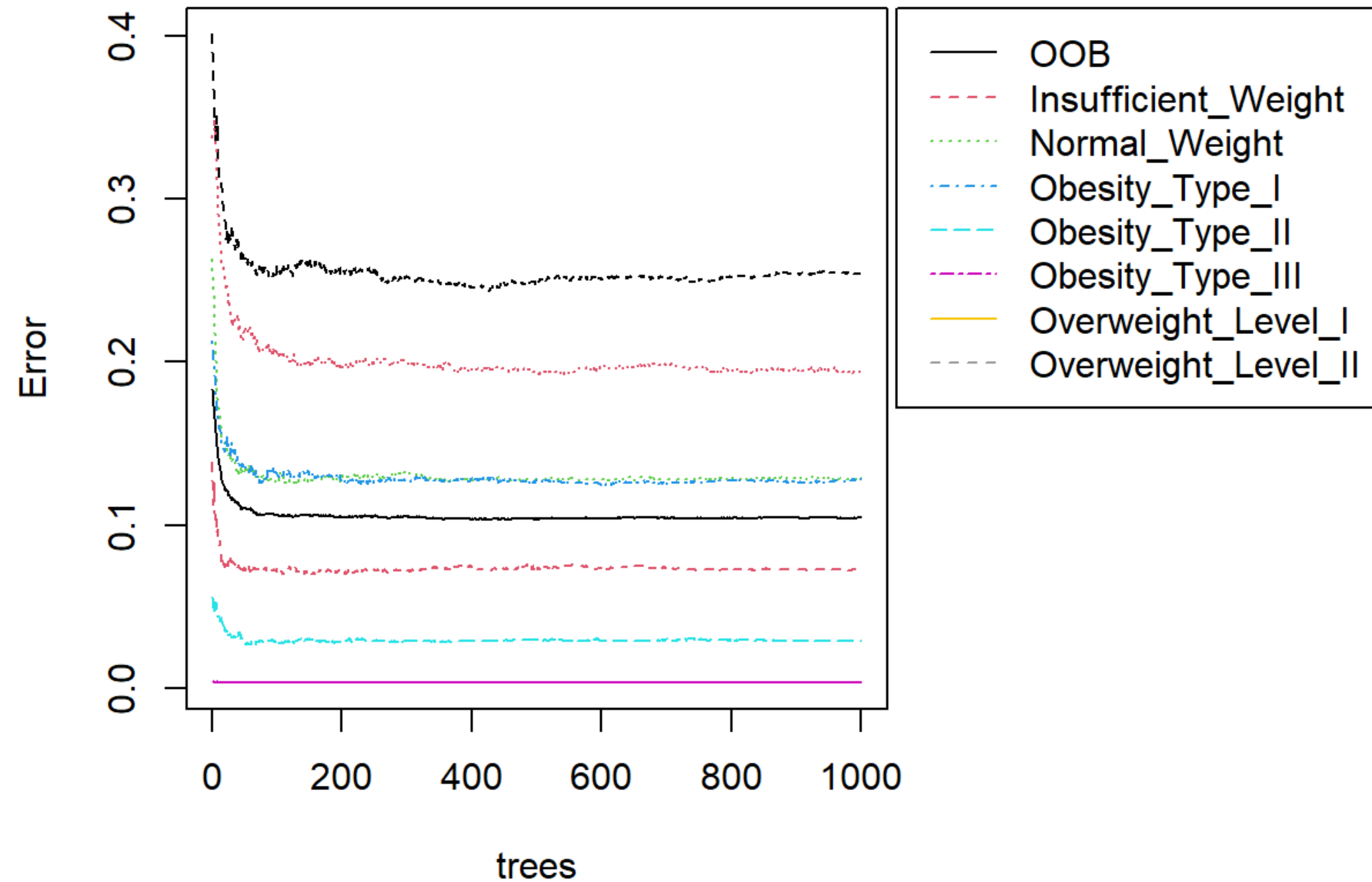
No Information Rate : 0.1949

P-Value [Acc > NIR] : < 2.2e-16

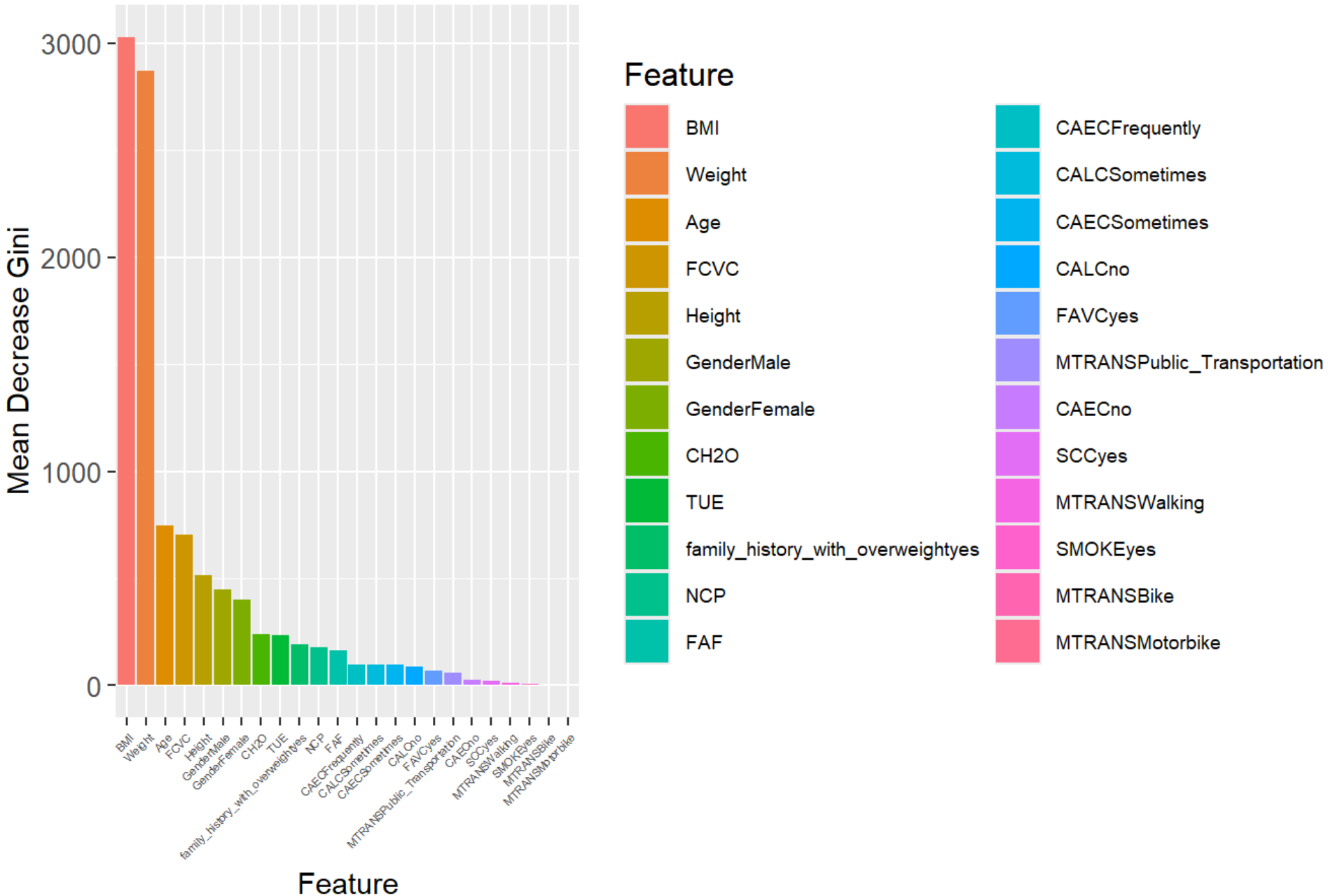
Kappa : 0.8875

Mcnemar's Test P-Value : NA

## Random Forest: Error per number of trees



# Feature Importances from Random Forest Model





# Improving the random forest model

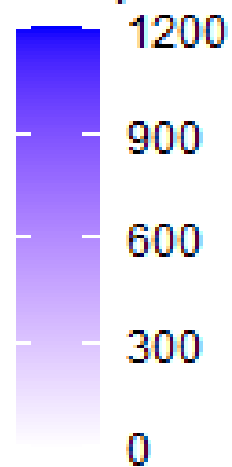
- Deleting features with low importance and high correlation

Confusion Matrix

True Label

Overweight_Level_II	0	15	52	7	0	61	621
Overweight_Level_I	6	74	16	0	0	550	82
Obesity_Type_III	0	0	1	0	1211	1	0
Obesity_Type_II	0	0	15	950	1	0	8
Obesity_Type_I	0	3	782	29	4	13	42
Normal_Weight	43	821	1	0	0	48	11
Insufficient_Weight	700	51	0	0	0	4	1

Freq



Accuracy: 0.9054

Insufficient\_Weight  
Normal\_Weight  
Obesity\_Type\_I  
Obesity\_Type\_II  
Obesity\_Type\_III  
Overweight\_Level\_I  
Overweight\_Level\_II

Comparision: Accuracy

