# Preparation

## **Data:**

1) TCGA-LIHC.htseq_fpkm.tsv
2) TCGA-LIHC.GDC_phenotype.tsv

```r
1   # Set the directory path
2   dir_project1 <- "C:/Users/choyo/Desktop/CODEDATA/biomeddatasci/II/hw2/"
3   dir_res1 <- dir_project1
4
5   # load in Data
6   f_data1 <- paste(dir_project1, 'TCGA-LIHC.htseq_fpkm.tsv', sep='/')
7   f_phe1 <- paste(dir_project1, 'TCGA-LIHC.GDC_phenotype.tsv', sep='/')
8   data1 <- read.delim(f_data1, sep='\t')
9   phe1 <- read.delim(f_phe1, sep='\t')
10
11  xt1 <- phe1[,99]
12
13  dim(data1)
14  colnames(data1)
15  dim(phe1)
16  colnames(phe1)
17
18  pro_id1 <- data1[,1] # probe_id
19  mode(pro_id1)
20  data1 <- data1[,-1]
21  mode(data1)
22  data1 <- as.matrix(data1) # convert data into numeric matrix
23  dim(data1)
24
25  r_sd1 <- apply(data1, 1, sd) # row (probe/gene) standard deviation (sd)
26  idx_row1 <- which(r_sd1 >= (sort(r_sd1, decreasing=T)[1000]))
27  data2 <- data1[idx_row1,] # 424 samples
28  dim(data2)
29  data2 <- t(data2)
30  dim(data2)
31
32  # define lables
33  # align the data with the phenotype data
34  id_data <- rownames(data2) # e.g., "TCGA.DD.A4NG.01A"
35  id_phe <- as.character(phe1[,1]) # e.g., "TCGA-DD-AAVQ-01A"
36  id_phe <- gsub('-', '.', id_phe, fixed=T) # change the '-' to '.' in the Phenothype sample id
37  n_t1 <- dim(data2)[1] # No. of samples
38  idx_order1 <- rep(0, n_t1)|
39- for (i in 1:n_t1){
40     idx_order1[i] <- which(id_phe %in% id_data[i])
41- }
42
43  #Prep Data
44
45  #sum(idx_order1 < 1) # make sure no sample (without phenotype information)
46  phe2 <- phe1[idx_order1,] # get the phenotype data of all data samples
47  table(phe2[,25])  # fibrosis
48  table(phe2[,116]) # sample type
49  table(phe2[,95])  # stage
50  idx_t1 <- which(phe2[,25] %in% "0 - No Fibrosis" | phe2[,25] %in% "1,2 - Portal Fibrosis")
51  idx_t2 <- which(phe2[,95] %in% "stage i" | phe2[,95] %in% "stage ii")
52  idx_T1 <- which(phe2[,116] %in% 'Primary Tumor') # tumor
53  idx_N1 <- which(phe2[,116] %in% 'Solid Tissue Normal') # normal
54  length(idx_T1)
55  length(idx_N1)
56  idx_C1 <- intersect(idx_t1, idx_t2) # cluster 1
57  idx_C2 <- setdiff(c(1:dim(data2)[1]), idx_C1) # clsuter 2
```

```
56   idx_C1 <- intersect(idx_t1, idx_t2) # cluster 1
57   idx_C2 <- setdiff(c(1:dim(data2)[1]), idx_C1) # clsuter 2
58   length(idx_C1)
59   length(idx_C2)
60   label3 <- rep('test', dim(data2)[1])
61   label3[idx_C1] <- 'Better_Fibrosis'
62   label3[idx_C2] <- 'Worse_Fibrosis'
63   label3 <- as.factor(label3)
64   # data <- as.data.frame(cbind(label3, data2))
65   data <- as.data.frame(data2)
66   data['label3'] <- label3
67   n_sample <- dim(data2)[1]
68   idx_train <- sample(c(1:n_sample), round(n_sample * 0.67))
69   idx_test <- setdiff(c(1:n_sample), idx_train)
70   data_train <- data[idx_train, ]
71   data_test <- data[idx_test,]
```
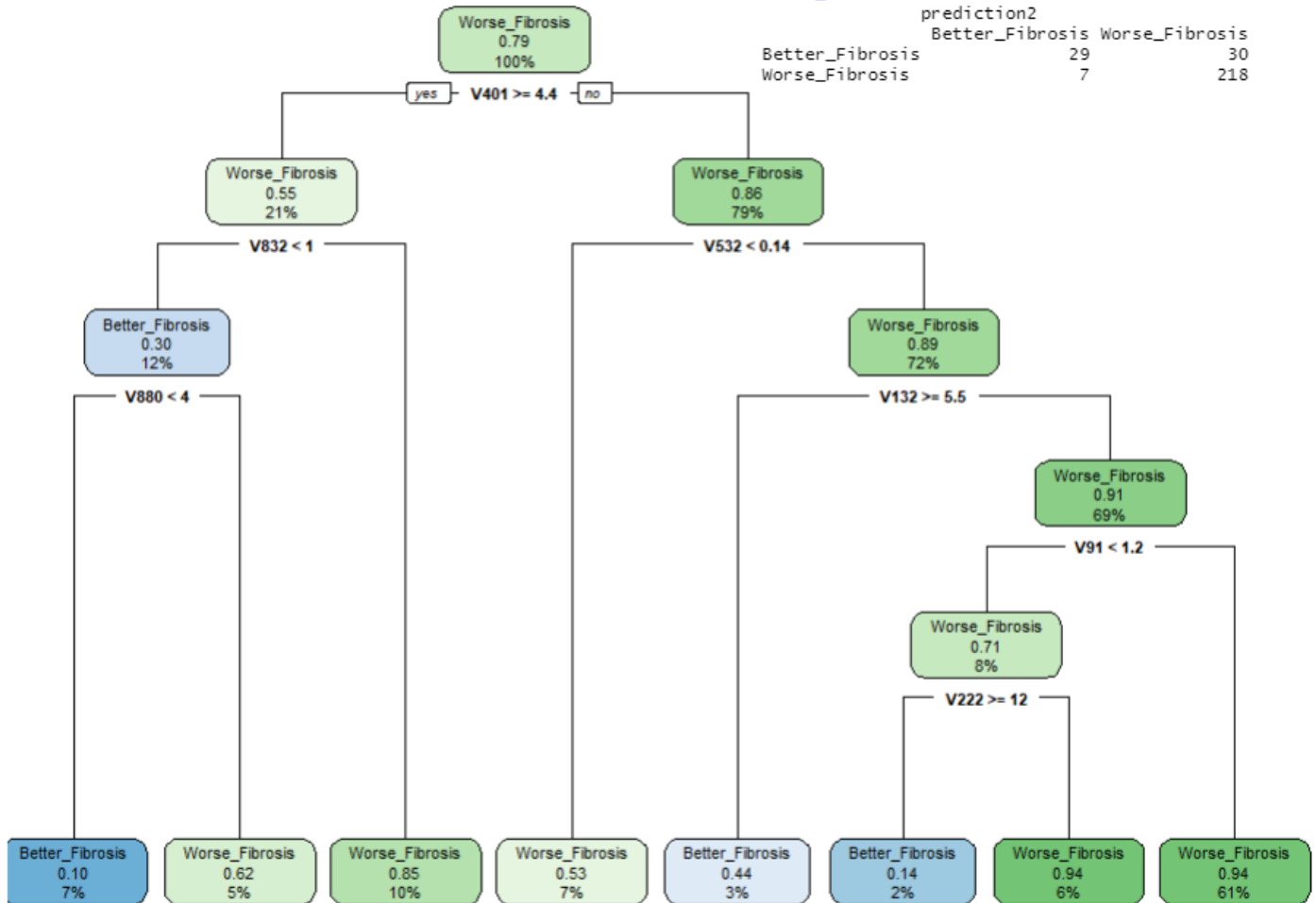
# Decision Tree

```
75   #Decision Tree
76   library(rpart)
77   library(rpart.plot)
78   n_sample <- dim(data2)[1]
79   idx_train <- sample(c(1:n_sample), round(n_sample * 0.67))
80   idx_test <- setdiff(c(1:n_sample), idx_train)
81   data_train <- data[idx_train, ]
82   data_test <- data[idx_test,]
83   # training
84   fit1 <- rpart(label3 ~., data=data_train, method = 'class')
85   rpart.plot(fit1)
86
87   # prediction
88   prediction1 <-predict(fit1, data_test, type = "class")
89
90
91   table_mat <- table(data_test$label3, prediction1)
92   table_mat
93   prediction2 <-predict(fit1, data_train, type = 'class')
94   table_mat <- table(data_train$label3, prediction2)
95   table_mat
```

```
> table_mat
                  prediction1
                  Better_Fibrosis Worse_Fibrosis
   Better_Fibrosis               3             29
   Worse_Fibrosis                9             99
> prediction2 <-predict(fit1, data_train, type = 'class')
> table_mat <- table(data_train$label3, prediction2)
> table_mat
                  prediction2
                  Better_Fibrosis Worse_Fibrosis
   Better_Fibrosis              29             30
   Worse_Fibrosis               7            218
```
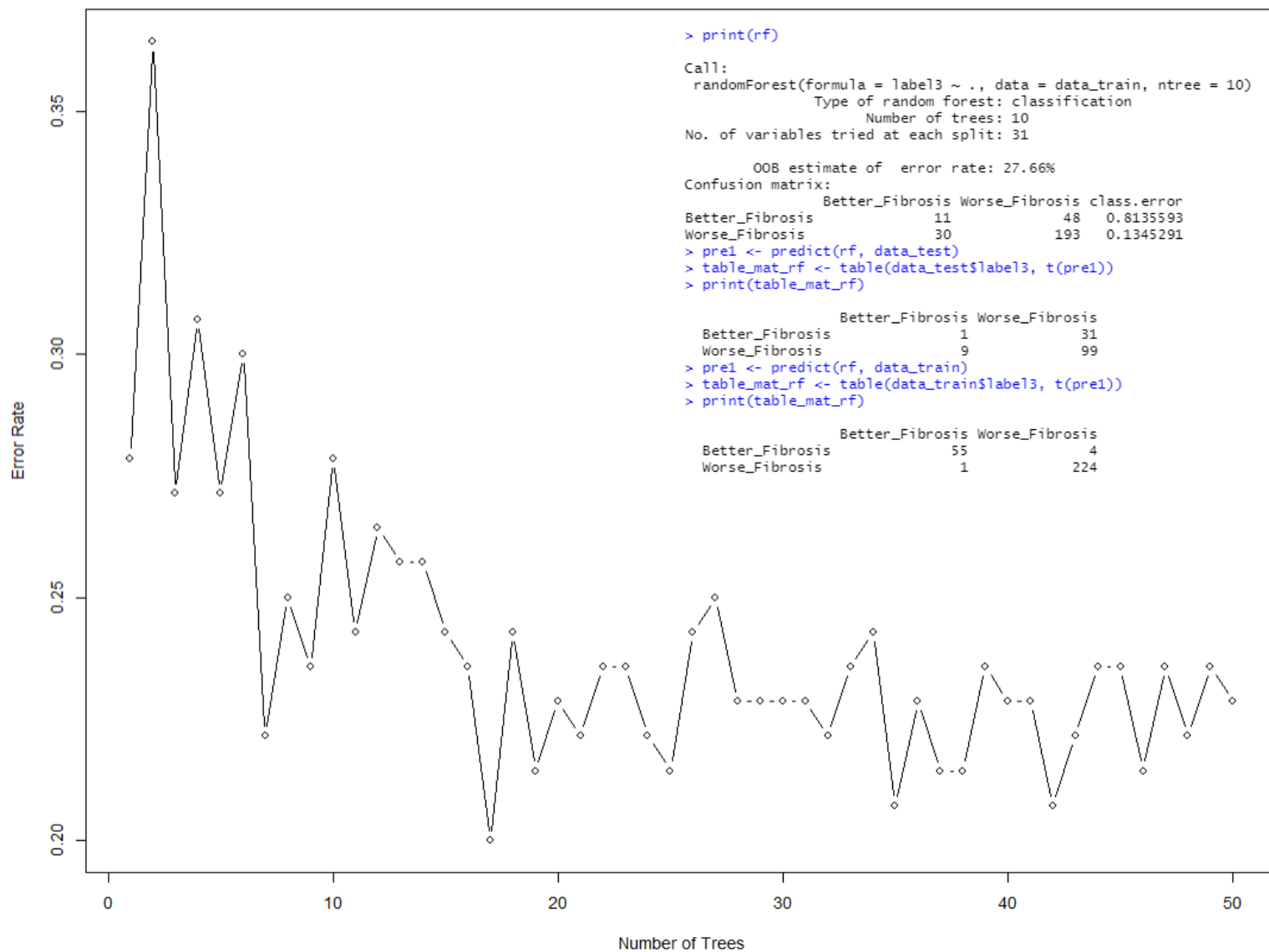
# Random Forest

```
 97  #Random forest
 98  #install.packages('randomForest')
 99  library('randomForest')
100  rf <- randomForest(label3~., data=data_train, ntree=10)
101  print(rf)
102  pre1 <- predict(rf, data_test)
103  table_mat_rf <- table(data_test$label3, t(pre1))
104  print(table_mat_rf)
105  pre1 <- predict(rf, data_train)
106  table_mat_rf <- table(data_train$label3, t(pre1))
107  print(table_mat_rf)
108
109  # Create a trees vs error plot
110  error_rate <- rep(0, 50) # record error rate for each model
111- for (i in 1:50) {
112    rf <- randomForest(label3~., data=data_train, ntree=i)
113    pre1 <- predict(rf, data_test)
114    error_rate[i] <- 1 - sum(diag(table(data_test$label3, pre1))) / sum(table(data_test$label3, pre1))
115- }
116
117  plot(1:50, error_rate, type='b', xlab='Number of Trees', ylab='Error Rate')
```

```
> print(rf)

Call:
 randomForest(formula = label3 ~ ., data = data_train, ntree = 10)
               Type of random forest: classification
                     Number of trees: 10
No. of variables tried at each split: 31

        OOB estimate of  error rate: 27.66%
Confusion matrix:
                Better_Fibrosis Worse_Fibrosis class.error
Better_Fibrosis              11             48   0.8135593
Worse_Fibrosis               30            193   0.1345291
> pre1 <- predict(rf, data_test)
> table_mat_rf <- table(data_test$label3, t(pre1))
> print(table_mat_rf)

                  Better_Fibrosis Worse_Fibrosis
  Better_Fibrosis               1             31
  Worse_Fibrosis                9             99
> pre1 <- predict(rf, data_train)
> table_mat_rf <- table(data_train$label3, t(pre1))
> print(table_mat_rf)

                  Better_Fibrosis Worse_Fibrosis
  Better_Fibrosis              55              4
  Worse_Fibrosis               1             224
```

# 3-cross folds validation

```r
121  # 3 fold validation
122  #install.packages('caret')
123  library(caret)
124  # define 3 folds for cross-validation
125  #divide dataset into 3 folds randomly
126  set.seed(123) # set a seed for reproducibility
127  n_sample <- dim(data2)[1]
128  folds <- sample(1:3, n_sample, replace = TRUE) # randomly divide the data into 3 folds
129
130  #train the model using the first and second fold
131  train_idx <- which(folds != 3) # select the indices of the first and second fold for training
132  data_train <- data2[train_idx, ]
133  label_train <- label3[train_idx]
134
135  #Test the model on the 3rd fold
136  test_idx <- which(folds == 3) # select the indices of the third fold for testing
137  data_test <- data2[test_idx, ]
138  label_test <- label3[test_idx]
139
140  #evaluate models on the test datasets with confusion matrix, precision, recall and accuracy
141  library(class)
142  predicted_labels <- knn(data_train, data_test, label_train, k = 5) # make predictions on the test data using the trained model
143  conf_mat <- table(predicted_labels, label_test) # compute confusion matrix
144  conf_mat
145  precision <- diag(conf_mat)/colSums(conf_mat) # compute precision
146  precision
147  recall <- diag(conf_mat)/rowSums(conf_mat) # compute recall
148  recall
149  accuracy <- sum(diag(conf_mat))/sum(conf_mat) # compute accuracy
150  accuracy
```

```
> conf_mat
                  label_test
predicted_labels   Better_Fibrosis Worse_Fibrosis
  Better_Fibrosis                1              9
  Worse_Fibrosis                26             95
> precision <- diag(conf_mat)/colSums(conf_mat) # compute precision
> precision
Better_Fibrosis  Worse_Fibrosis
     0.03703704      0.91346154
> recall <- diag(conf_mat)/rowSums(conf_mat) # compute recall
> recall
Better_Fibrosis  Worse_Fibrosis
       0.100000        0.785124
> accuracy <- sum(diag(conf_mat))/sum(conf_mat) # compute accuracy
> accuracy
[1] 0.7328244
```