

## PSTAT 131 Homework 3

HaozeZhu &amp; Yubo Wei

05/21/22

```
##loading library
```

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.3      v purrr   0.3.4
## v tibble  3.1.1      v dplyr  1.0.5
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.5.1
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(ROCR)
library(tree)
```

```
## Registered S3 method overwritten by 'tree':
##   method      from
##   print.tree cli
```

```
library(maptree)
```

```
## Loading required package: cluster
```

```
## Loading required package: rpart
```

```
library(class)
library(lattice)
library(ggthemes)
library(superheat)
```

```
## loading the drug dataset and transform variables
```

```
drug_use <- read_csv('drug.csv',
col_names = c('ID', 'Age', 'Gender', 'Education', 'Country', 'Ethnicity', 'Nscore', 'Escore', 'Oscore', 'Ascore',
'Choc', 'Coke', 'Crack', 'Ecstasy', 'Heroin', 'Ketamine', 'Legalh', 'LSD', 'Meth', 'Mushrooms', 'Nicotine', 'Seme'
```

```
##
## -- Column specification -----
## cols(
##   .default = col_character(),
##   ID = col_double(),
##   Age = col_double(),
##   Gender = col_double(),
##   Education = col_double(),
##   Country = col_double(),
##   Ethnicity = col_double(),
##   Nscore = col_double(),
##   Escore = col_double(),
##   Oscore = col_double(),
##   Ascore = col_double(),
##   Cscore = col_double(),
##   Impulsive = col_double(),
##   SS = col_double()
## )
## i Use 'spec()' for the full column specifications.
```

```
drug_use <- drug_use %>% mutate_at(as.ordered, .vars=vars(Alcohol:VSA))
drug_use <- drug_use %>%
mutate(Gender = factor(Gender, labels=c("Male", "Female"))) %>% mutate(Ethnicity = factor(Ethnicity, labels=c("White", "Black", "Hispanic", "Asian", "Other")))
```

##1. Logistic regression for drug use prediction ###a

```
## create a new factor response variable recent_cannabis_use
drug_use = drug_use %>% mutate(recent_cannabis_use=ifelse(Cannabis >= "CL3", "Yes", "No")) %>% mutate(recent_cannabis_use = factor(recent_cannabis_use, labels=c("Yes", "No")))
```

###b

```
##create a new tibble
drug_use_subset <- drug_use %>% select(Age:SS, recent_cannabis_use)
##Split drug_use_subset into a training data set and a test data set
set.seed(1)
training.indices = sample(1:nrow(drug_use_subset), 1500)
drug_use_train = drug_use_subset[training.indices,]
drug_use_test = drug_use_subset[-training.indices,]
dim(drug_use_train)
```

```
## [1] 1500 13
```

```
dim(drug_use_test)
```

```
## [1] 385 13
```

###c

```
##Fit a logistic regression
glm.fit = glm(recent_cannabis_use ~ ., data=drug_use_train, family=binomial)
# Summarize the logistic regression model
summary(glm.fit)
```

```
##
## Call:
## glm(formula = recent_cannabis_use ~ ., family = binomial, data = drug_use_train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.9072  -0.5971   0.1416   0.5426   2.6600
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.94949    0.64574   1.470 0.141457
## Age           -0.84406    0.09328  -9.049 < 2e-16 ***
## GenderFemale  -0.55929    0.15715  -3.559 0.000372 ***
## Education     -0.33389    0.07962  -4.193 2.75e-05 ***
## CountryCanada 13.10904   627.22755   0.021 0.983325
## CountryNew Zealand -1.16844    0.31848  -3.669 0.000244 ***
## CountryOther  -0.05676    0.46772  -0.121 0.903412
## CountryIreland -0.28763    0.67573  -0.426 0.670354
## CountryUK     -0.43371    0.37043  -1.171 0.241674
## CountryUSA    -1.75636    0.19262  -9.118 < 2e-16 ***
## EthnicityAsian -0.67025    0.96037  -0.698 0.485230
## EthnicityWhite  0.74053    0.63843   1.160 0.246081
## EthnicityMixed:White/Black -0.04713    1.09013  -0.043 0.965515
## EthnicityOther  1.07889    0.76823   1.404 0.160206
## EthnicityMixed:White/Asian  0.72525    1.01565   0.714 0.475178
## EthnicityMixed:Black/Asian 14.27149   766.28165   0.019 0.985141
## Nscore        -0.10143    0.09034  -1.123 0.261551
## Escore        -0.13375    0.09559  -1.399 0.161742
## Oscore         0.71000    0.09137   7.770 7.83e-15 ***
## Ascore         0.03058    0.08232   0.372 0.710251
## Cscore        -0.35855    0.09132  -3.926 8.63e-05 ***
## Impulsive     -0.09043    0.10093  -0.896 0.370290
## SS            0.58068    0.10836   5.359 8.39e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 2077.2  on 1499  degrees of freedom
## Residual deviance: 1202.1  on 1477  degrees of freedom
## AIC: 1248.1
##
## Number of Fisher Scoring iterations: 14
```

##2. Decision tree models of drug use

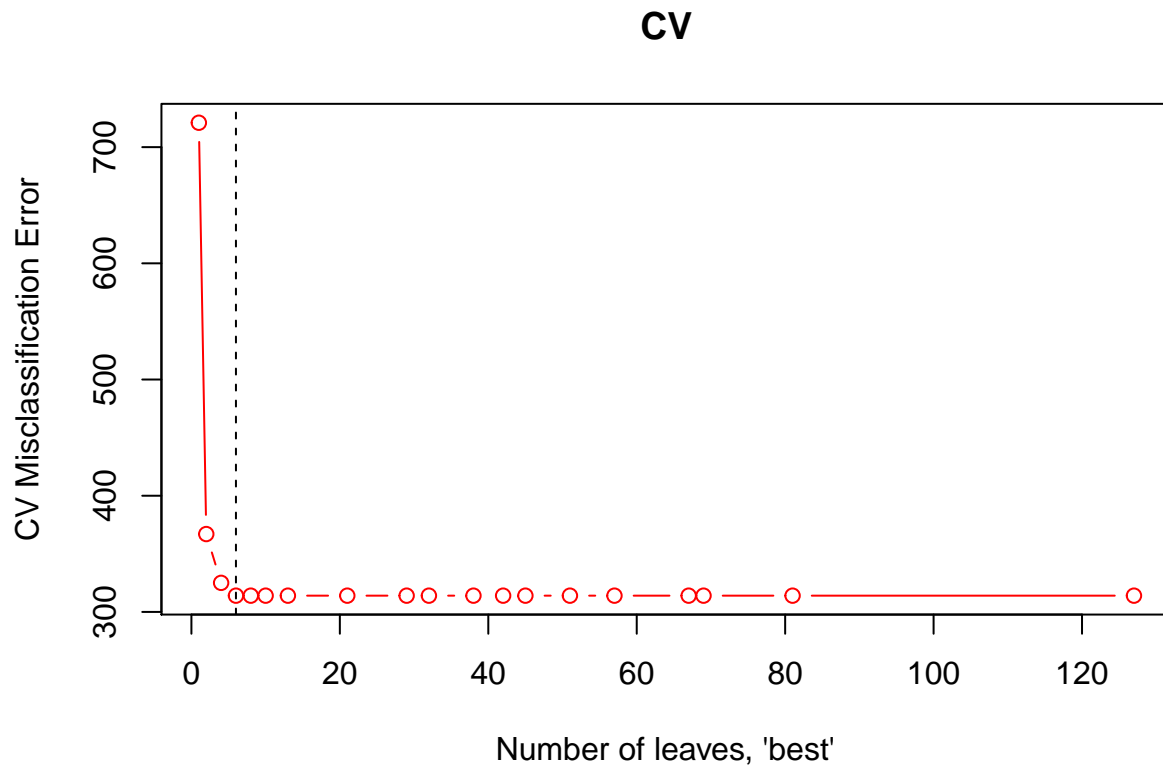
```
## Construct a decision tree
tree_parameters = tree.control(nobs=nrow(drug_use_train), minsize=10, mindev=1e-3)
drugtree = tree(recent_cannabis_use~., control = tree_parameters, data = drug_use_train)
```

###a Use 10-fold CV to select the a tree

```

# Set random seed
set.seed(1)
cv = cv.tree(drugtree, FUN=prune.misclass, K=10)
best.size.cv = 6
plot(cv$size, cv$dev, type="b",
     xlab = "Number of leaves, \'best\'", ylab = "CV Misclassification Error",
     col = "red", main="CV")
abline(v=best.size.cv, lty=2)

```

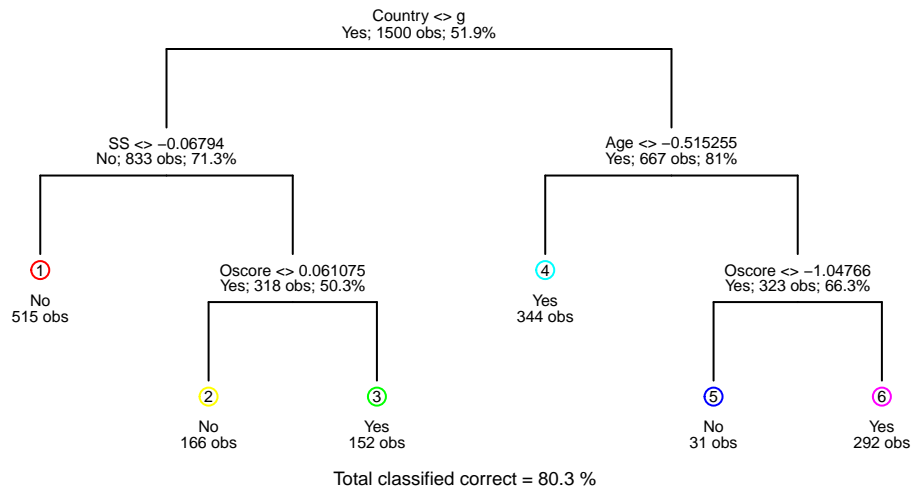


```
###b Prune the tree
```

```

# prune the original tree using the best size in a
drugtree.pruned = prune.misclass(drugtree, best=best.size.cv)
draw.tree(drugtree.pruned, nodeinfo=TRUE, cex = 0.5)

```



The variable 'Country' is split

first in this decision tree.

###c Compute and print the confusion matrix for the test data

```

# Predict on test set
predictions = predict(drugtree.pruned, drug_use_test, type="class")
# get the true response of the test data
truth = drug_use_test$recent_cannabis_use
# Obtain confusion matrix
confusion_matrix = table(truth, predictions)
confusion_matrix

```

```

##      predictions
## truth  No Yes
##   No  125  40
##   Yes   45 175

```

```

# get true positive rate
true_positive_rate = confusion_matrix[2,2]/sum(confusion_matrix[2,])
true_positive_rate

```

```
## [1] 0.7954545
```

```

# get false positive rate
false_positive_rate = confusion_matrix[1,2]/sum(confusion_matrix[1,])
false_positive_rate

```

```
## [1] 0.2424242
```

###3.Model Comparison ###a

```

# get prediction from logistic regression model using test data
prob_log_testing = predict(glm.fit,drug_use_test,type="response")
pred_log = prediction(prob_log_testing, truth)
#calculate the True Positive Rate and False Positive Rate by performance()
perf_log = performance(pred_log, measure="tpr", x.measure="fpr")
# for the decision tree model

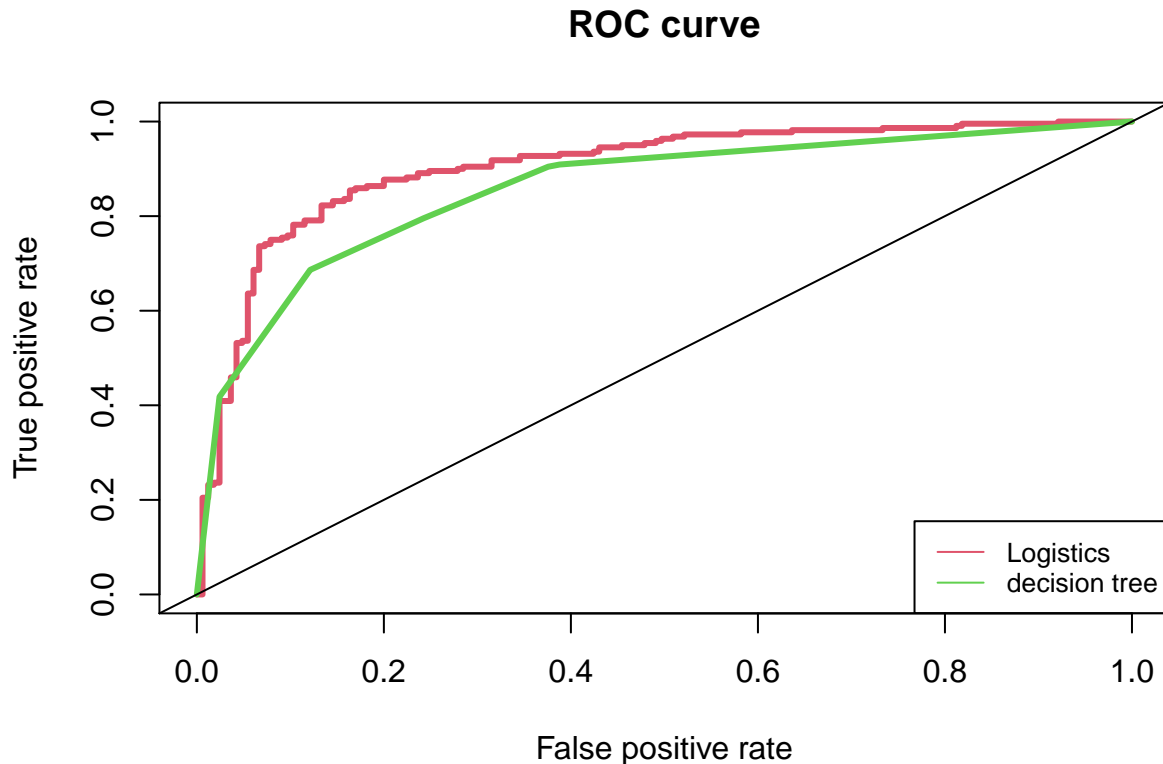
```

```

prob_tree_testing = predict(drugtree.pruned,newdata=drug_use_test)[,2]
pred_tree = prediction(prob_tree_testing,truth)
perf_tree = performance(pred_tree, measure="tpr", x.measure="fpr")

plot(perf_log, col=2, lwd=3, main="ROC curve")
plot(perf_tree, col=3, lwd=3,add=TRUE)
legend(x="bottomright",legend=c("Logistics", "decision tree"),
      col=c(2, 3), lty=1, cex=0.8)
abline(0,1)

```



```
###b
```

```

##Compute the AUC for logistic regression model
auc_log = performance(pred_log, "auc")@y.values
auc_log

```

```

## [[1]]
## [1] 0.902562

```

```

##Compute the AUC for decision tree model
auc_tree= performance(pred_tree, "auc")@y.values
auc_tree

```

```

## [[1]]
## [1] 0.8570523

```

The logistic regression model has larger AUC, which suggests overall the logistic regression model is a better model.

##4. Clustering and dimension reduction for gene expression data

```
## read leukemia_data
leukemia_data <- read_csv("leukemia_data.csv")

## Warning: Duplicated column names deduplicated: 'FCGRT' => 'FCGRT_1' [3],
## 'TUBB4B' => 'TUBB4B_1' [49], 'SSR1' => 'SSR1_1' [67], 'HSP90AB1' =>
## 'HSP90AB1_1' [115], 'TMBIM6' => 'TMBIM6_1' [118], 'GAB1' => 'GAB1_1' [119],
## 'MPHOSPH9' => 'MPHOSPH9_1' [153], 'STK38' => 'STK38_1' [157], 'SFPQ' =>
## 'SFPQ_1' [159], 'RIPOR2' => 'RIPOR2_1' [181], 'HLA-F' => 'HLA-F_1' [188],
## 'PRPF40A' => 'PRPF40A_1' [198], 'SEPT6' => 'SEPT6_1' [205], 'CD22' =>
## 'CD22_1' [235], 'NCF4' => 'NCF4_1' [250], 'WAS' => 'WAS_1' [260], 'HLA-
## G' => 'HLA-G_1' [297], 'TRAF3IP3' => 'TRAF3IP3_1' [307], 'ZNF266' =>
## 'ZNF266_1' [364], 'CRYBG1' => 'CRYBG1_1' [441], 'BRD8' => 'BRD8_1' [460], 'MDC1'
## => 'MDC1_1' [464], 'RAC2' => 'RAC2_1' [478], 'IL10RB' => 'IL10RB_1' [483],
## 'AKAP17A' => 'AKAP17A_1' [542], 'N4BP2L1' => 'N4BP2L1_1' [547], 'ARPC4' =>
## 'ARPC4_1' [565], 'SRSF10' => 'SRSF10_1' [576], 'RAPGEF2' => 'RAPGEF2_1' [583],
## 'PARP2' => 'PARP2_1' [587], 'TRIM33' => 'TRIM33_1' [610], 'KAT8' =>
## 'KAT8_1' [665], 'ASMTL' => 'ASMTL_1' [715], 'LSM7' => 'LSM7_1' [727],
## 'HLA-DQB1' => 'HLA-DQB1_1' [732], 'FMR1' => 'FMR1_1' [826], 'RASGRP2' =>
## 'RASGRP2_1' [858], 'LIMK2' => 'LIMK2_1' [866], 'TMEM106C' => 'TMEM106C_1' [881],
## 'TGOLN2' => 'TGOLN2_1' [937], 'SLC25A1' => 'SLC25A1_1' [940], 'NMT1' =>
## 'NMT1_1' [942], 'ENSA' => 'ENSA_1' [947], 'ENSA' => 'ENSA_2' [948], 'UBR5'
## => 'UBR5_1' [963], 'UBE2J1' => 'UBE2J1_1' [966], 'ACTN1' => 'ACTN1_1' [994],
## 'TRA2A' => 'TRA2A_1' [1003], 'ATXN10' => 'ATXN10_1' [1057], 'CUL1' =>
## 'CUL1_1' [1077], 'XBP1' => 'XBP1_1' [1094], 'ATP2A2' => 'ATP2A2_1' [1110],
## 'LDLRAD4' => 'LDLRAD4_1' [1118], 'ARHGEF2' => 'ARHGEF2_1' [1134],
## 'IDH3B' => 'IDH3B_1' [1141], 'SERBP1' => 'SERBP1_1' [1188], 'TRIM44' =>
## 'TRIM44_1' [1205], 'TRIM44' => 'TRIM44_2' [1206], 'PTPRC' => 'PTPRC_1' [1219],
## 'PTPRC' => 'PTPRC_2' [1220], 'PPP2R5C' => 'PPP2R5C_1' [1235], 'PPP2R5C'
## => 'PPP2R5C_2' [1236], 'ADAM10' => 'ADAM10_1' [1241], 'NFATC3' =>
## 'NFATC3_1' [1252], 'ILF3' => 'ILF3_1' [1264], 'RBM6' => 'RBM6_1' [1274],
## 'CTNNA1' => 'CTNNA1_1' [1297], 'CTNNA1' => 'CTNNA1_2' [1298], 'IGHM' =>
## 'IGHM_1' [1302], 'IGHM' => 'IGHM_2' [1303], 'IGHM' => 'IGHM_3' [1304], 'SFPQ' =>
## 'SFPQ_2' [1321], 'RBCK1' => 'RBCK1_1' [1398], 'NFATC2IP' => 'NFATC2IP_1' [1408],
## 'ILF3' => 'ILF3_2' [1432], 'RAE1' => 'RAE1_1' [1436], 'ITPR1' =>
## 'ITPR1_1' [1443], 'NCBP2' => 'NCBP2_1' [1448], 'STAT1' => 'STAT1_1' [1486],
## 'AZIN1' => 'AZIN1_1' [1497], 'SEC13' => 'SEC13_1' [1517], 'ABI1' =>
## 'ABI1_1' [1565], 'CYB5B' => 'CYB5B_1' [1607], 'HUWE1' => 'HUWE1_1' [1624],
## 'RAB1A' => 'RAB1A_1' [1634], 'AHCYL1' => 'AHCYL1_1' [1652], 'EIF1AX' =>
## 'EIF1AX_1' [1661], 'MAGED2' => 'MAGED2_1' [1689], 'SCAF11' => 'SCAF11_1' [1709],
## 'BLCAP' => 'BLCAP_1' [1716], 'TROVE2' => 'TROVE2_1' [1729], 'CTCF' =>
## 'CTCF_1' [1745], 'RAB8A' => 'RAB8A_1' [1754], 'ACTR2' => 'ACTR2_1' [1768],
## 'HMGN4' => 'HMGN4_1' [1771], 'NDUFB7' => 'NDUFB7_1' [1793], 'VAMP3' =>
## 'VAMP3_1' [1796], 'SRSF6' => 'SRSF6_1' [1808], 'TNPO3' => 'TNPO3_1' [1811],
## 'SRSF1' => 'SRSF1_1' [1834], 'TMED10' => 'TMED10_1' [1847], 'AP3D1' =>
## 'AP3D1_1' [1872], 'MAPKAPK2' => 'MAPKAPK2_1' [1877], 'BRD2' => 'BRD2_1' [1891],
## 'BRD2' => 'BRD2_2' [1892], 'GARS' => 'GARS_1' [1901], 'SNX1' => 'SNX1_1' [1902],
## 'TSC22D3' => 'TSC22D3_1' [1927], 'AMD1' => 'AMD1_1' [1951], 'LITAF' =>
## 'LITAF_1' [2011], 'GLUD1' => 'GLUD1_1' [2059], 'KDELRL1' => 'KDELRL1_1' [2079],
## 'PGK1' => 'PGK1_1' [2099], 'VDAC2' => 'VDAC2_1' [2107], 'ADH5' =>
## 'ADH5_1' [2111], 'MEF2C' => 'MEF2C_1' [2113], 'MEF2C' => 'MEF2C_2' [2114],
## 'RCN2' => 'RCN2_1' [2125], 'PCMT1' => 'PCMT1_1' [2134], 'PCMT1' =>
## 'PCMT1_2' [2135], 'CD79A' => 'CD79A_1' [2149], 'MARCH6' => 'MARCH6_1' [2169],
## 'CBX3' => 'CBX3_1' [2180], 'LSM14A' => 'LSM14A_1' [2217], 'SORL1' =>
```

```

## 'SORL1_1' [2220], 'ICAM2' => 'ICAM2_1' [2244], 'SNRPB' => 'SNRPB_1' [2246],
## 'CYB5A' => 'CYB5A_1' [2248], 'BTN3A2' => 'BTN3A2_1' [2277], 'DICER1' =>
## 'DICER1_1' [2280], 'HADH' => 'HADH_1' [2281], 'HDGF' => 'HDGF_1' [2285], 'SEPT6'
## => 'SEPT6_2' [2306], 'SSBP1' => 'SSBP1_1' [2315], 'H2AFV' => 'H2AFV_1' [2318],
## 'PTPA' => 'PTPA_1' [2331], 'FBL' => 'FBL_1' [2354], 'OGT' => 'OGT_1' [2362],
## 'SLC25A1' => 'SLC25A1_2' [2377], 'FUBP1' => 'FUBP1_1' [2386], 'TUBGCP2' =>
## 'TUBGCP2_1' [2400], 'COX5B' => 'COX5B_1' [2402], 'VDAC1' => 'VDAC1_1' [2410],
## 'HNRNPDL' => 'HNRNPDL_1' [2431], 'THUMP1' => 'THUMP1_1' [2443], 'CDV3'
## => 'CDV3_1' [2444], 'UBE3B' => 'UBE3B_1' [2447], 'SFPQ' => 'SFPQ_3' [2451],
## 'STX16' => 'STX16_1' [2452], 'SMARCA2' => 'SMARCA2_1' [2471], 'CHD8' =>
## 'CHD8_1' [2475], 'TCF25' => 'TCF25_1' [2490], 'API5' => 'API5_1' [2491],
## 'SAP18' => 'SAP18_1' [2493], 'AHCYL1' => 'AHCYL1_2' [2501], 'CTBP1' =>
## 'CTBP1_1' [2503], 'AES' => 'AES_1' [2512], 'PURA' => 'PURA_1' [2514], 'BCL11A'
## => 'BCL11A_1' [2518], 'BUB3' => 'BUB3_1' [2534], 'RER1' => 'RER1_1' [2537],
## 'ATXN2L' => 'ATXN2L_1' [2541], 'JAK1' => 'JAK1_1' [2548], 'GUSBP11' =>
## 'GUSBP11_1' [2564], 'JTB' => 'JTB_1' [2568], 'BRD3' => 'BRD3_1' [2571], 'RSU1'
## => 'RSU1_1' [2584], 'ADD3' => 'ADD3_1' [2619], 'UBE2I' => 'UBE2I_1' [2627],
## 'MRPS12' => 'MRPS12_1' [2640], 'CTNNA1' => 'CTNNA1_3' [2641], 'XRCC5' =>
## 'XRCC5_1' [2642], 'ITGA4' => 'ITGA4_1' [2644], 'CTNNA1' => 'CTNNA1_4' [2647],
## 'FYN' => 'FYN_1' [2649], 'ERG' => 'ERG_1' [2652], 'RAC1' => 'RAC1_1' [2654],
## 'LCK' => 'LCK_1' [2657], 'PTK2B' => 'PTK2B_1' [2664], 'SKP1' =>
## 'SKP1_1' [2665], 'PRKDC' => 'PRKDC_1' [2666], 'MYC' => 'MYC_1' [2668], 'RBL2'
## => 'RBL2_1' [2673], 'AZIN1' => 'AZIN1_2' [2674], 'CCNA2' => 'CCNA2_1' [2681],
## 'FOS' => 'FOS_1' [2688], 'FOS' => 'FOS_2' [2689], 'RAF1' => 'RAF1_1' [2690],
## 'RAP1B' => 'RAP1B_1' [2692], 'ERCC1' => 'ERCC1_1' [2696], 'ERCC1' =>
## 'ERCC1_2' [2697], 'RAN' => 'RAN_1' [2702], 'TRIM27' => 'TRIM27_1' [2703],
## 'PMS2P3' => 'PMS2P3_1' [2708], 'TGFB2' => 'TGFB2_1' [2710], 'PCNA' =>
## 'PCNA_1' [2712], 'MYC' => 'MYC_2' [2714], 'CDK13' => 'CDK13_1' [2717],
## 'CCND3' => 'CCND3_1' [2719], 'FARSA' => 'FARSA_1' [2732], 'FARSA' =>
## 'FARSA_2' [2733], 'DAXX' => 'DAXX_1' [2734], 'UBE3A' => 'UBE3A_1' [2735],
## 'ARAF' => 'ARAF_1' [2739], 'UBE2N' => 'UBE2N_1' [2747], 'RASA1' =>
## 'RASA1_1' [2748], 'ABL1' => 'ABL1_1' [2749], 'ABL1' => 'ABL1_2' [2750], 'MTA1'
## => 'MTA1_1' [2753], 'EIF3I' => 'EIF3I_1' [2754], 'SYK' => 'SYK_1' [2761],
## 'TOP2A' => 'TOP2A_1' [2762], 'RB1' => 'RB1_1' [2764], 'TOP2B' =>
## 'TOP2B_1' [2765], 'TNFRSF1B' => 'TNFRSF1B_1' [2766], 'GRB2' => 'GRB2_1' [2769],
## 'RBM5' => 'RBM5_1' [2770], 'N4BP2L1' => 'N4BP2L1_2' [2773], 'N4BP2L2' =>
## 'N4BP2L2_1' [2774], 'NME1' => 'NME1_1' [2775], 'TYMS' => 'TYMS_1' [2776],
## 'DYRK1A' => 'DYRK1A_1' [2778], 'FEN1' => 'FEN1_1' [2779], 'FEN1' =>
## 'FEN1_2' [2780], 'ETS2' => 'ETS2_1' [2781], 'FNTA' => 'FNTA_1' [2783], 'JAK1'
## => 'JAK1_2' [2787], 'MYB' => 'MYB_1' [2792], 'MYB' => 'MYB_2' [2793], 'MYB' =>
## 'MYB_3' [2794], 'MYB' => 'MYB_4' [2795], 'MYB' => 'MYB_5' [2796], 'SMAD2' =>
## 'SMAD2_1' [2798], 'PTEN' => 'PTEN_1' [2799], 'MAPKAPK2' => 'MAPKAPK2_2' [2800],
## 'PSMD9' => 'PSMD9_1' [2801], 'PSMA4' => 'PSMA4_1' [2806], 'SRF' =>
## 'SRF_1' [2810], 'LYN' => 'LYN_1' [2815], 'IL7R' => 'IL7R_1' [2817], 'TCF3' =>
## 'TCF3_1' [2818], 'TCF3' => 'TCF3_2' [2819], 'NFKB1' => 'NFKB1_1' [2820], 'NFKB1'
## => 'NFKB1_2' [2821], 'RPA1' => 'RPA1_1' [2822], 'PPP2R2A' => 'PPP2R2A_1' [2823],
## 'TERF1' => 'TERF1_1' [2826], 'BCR' => 'BCR_1' [2828], 'RBBP4' =>
## 'RBBP4_1' [2830], 'TERF2' => 'TERF2_1' [2831], 'PSMB4' => 'PSMB4_1' [2834],
## 'PSMB7' => 'PSMB7_1' [2836], 'PARP1' => 'PARP1_1' [2838], 'RELA' =>
## 'RELA_1' [2840], 'RELA' => 'RELA_2' [2841], 'EIF2S3' => 'EIF2S3_1' [2842],
## 'YWHAZ' => 'YWHAZ_1' [2846], 'PTP4A2' => 'PTP4A2_1' [2847], 'POLR2H' =>
## 'POLR2H_1' [2850], 'GAB1' => 'GAB1_2' [2851], 'PRKDC' => 'PRKDC_2' [2852],
## 'PRKCB' => 'PRKCB_1' [2855], 'SAT1' => 'SAT1_1' [2862], 'PTPRE' =>
## 'PTPRE_1' [2865], 'RPL22' => 'RPL22_1' [2866], 'EIF2S1' => 'EIF2S1_1' [2867],

```



```
## 'CYC1' => 'CYC1_1' [2869], 'HSP90AB1' => 'HSP90AB1_2' [2870], 'CD44' =>
## 'CD44_1' [2873], 'MAP2K1' => 'MAP2K1_1' [2875], 'TNK2' => 'TNK2_1' [2877],
## 'GNA13' => 'GNA13_1' [2879], 'NR3C1' => 'NR3C1_1' [2882], 'RAB1A' =>
## 'RAB1A_2' [2888], 'ODC1' => 'ODC1_1' [2890], 'PLCG2' => 'PLCG2_1' [2891], 'RFC4'
## => 'RFC4_1' [2894], 'FLT3' => 'FLT3_1' [2895], 'EIF2AK2' => 'EIF2AK2_1' [2902],
## 'USP9X' => 'USP9X_1' [2913], 'PSMD7' => 'PSMD7_1' [2917], 'PPP1CA' =>
## 'PPP1CA_1' [2924], 'TUBB4B' => 'TUBB4B_2' [2926], 'ARRB2' => 'ARRB
```

```
##
## -- Column specification -----
## cols(
##   .default = col_double(),
##   Type = col_character()
## )
## i Use 'spec()' for the full column specifications.
```

```
####a
```

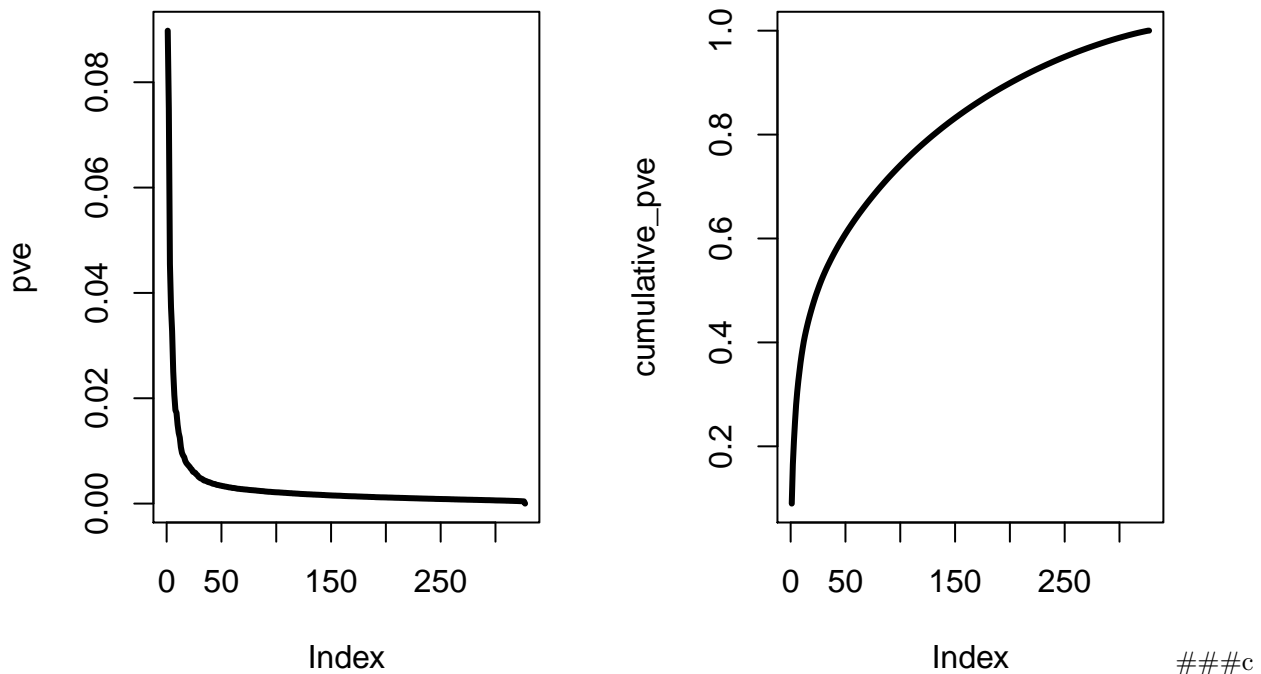
```
##Convert the Type column to factor
leukemia_data = leukemia_data %>% mutate(Type = factor(Type))
##the number of patients with each leukemia subtype
table(leukemia_data$Type)
```

```
##
##      BCR-ABL      E2A-PBX1 Hyperdip50      MLL      OTHERS      T-ALL      TEL-AML1
##           15           27           64           20           79           43           79
```

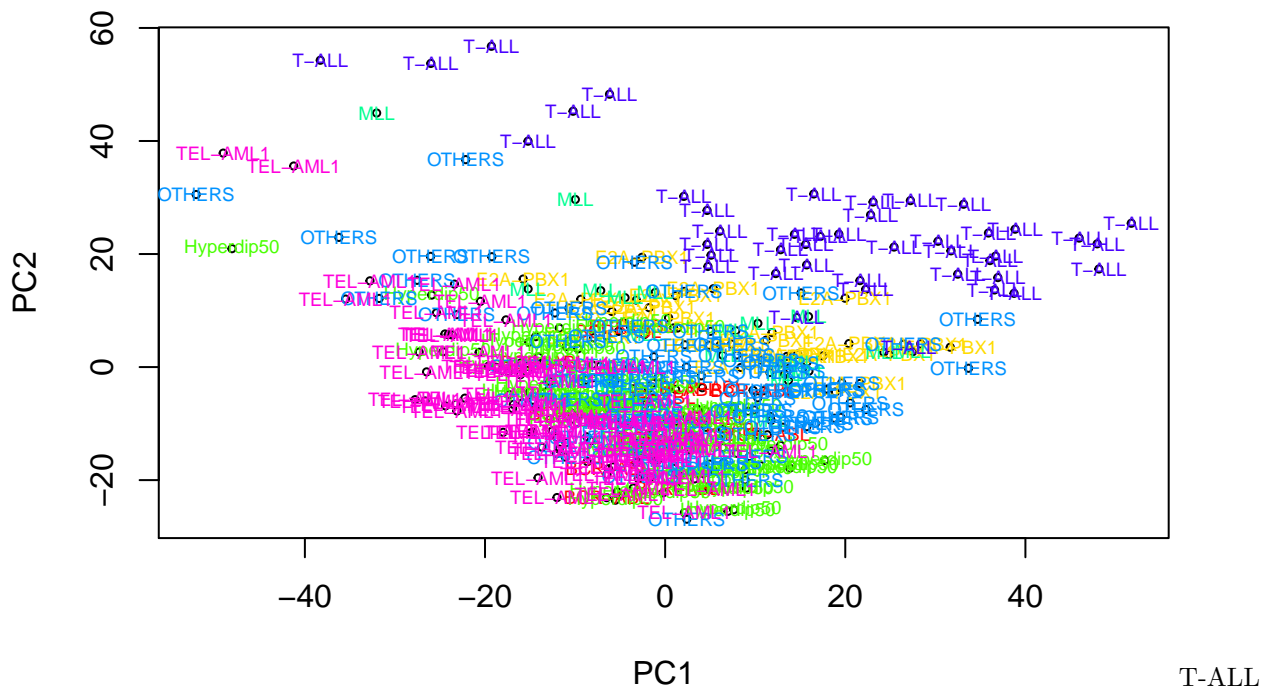
Subtype “BCR-ABL” occurs the least in this data

```
####b
```

```
## exclude the Type column
leukemia_data_wt = leukemia_data %>% select(-Type)
pr.out=prcomp(leukemia_data_wt, scale=TRUE, center = TRUE)
## get the variance explained by each principal component
pr.var=pr.out$sdev ^2
pve=pr.var/sum(pr.var)
## get cumulative pve
cumulative_pve = cumsum(pve)
## plot PVE and cumulative pve
par(mfrow=c(1, 2))
plot(pve, type="l", lwd=3)
plot(cumulative_pve, type="l", lwd=3)
```



```
## set color
rainbow_colors <- rainbow(7)
plot_colors <- rainbow_colors[leukemia_data$Type]
plot(pr.out$x[,1],pr.out$x[,2],xlab="PC1", ylab="PC2",cex=0.5)
text(pr.out$x[,1],pr.out$x[,2],labels = leukemia_data$Type,col=plot_colors,cex=0.6)
```



is most clearly separated from the others along the PC1 axis

```
## select PC1 value for loadings
pc_1 = pr.out$rotation[,1]
## sort the PC1 value
```

```
pc_1_sorted= sort(abs(pc_1),decreasing = TRUE)
head(pc_1_sorted,1)
```

```
##      SEMA3F
## 0.04517148
```

SEMA3F has the highest absolute loadings for PC1

```
## Print the first 6 genes in this sorted vector
head(pc_1_sorted,6)
```

```
##      SEMA3F      CCT2      LDHB      COX6C      SNRPD2      ELK3
## 0.04517148 0.04323818 0.04231619 0.04183480 0.04179822 0.04155821
```

```
###f
```

```
## load dendextend library
library(dendextend)
```

```
##
## -----
## Welcome to dendextend version 1.15.1
## Type citation('dendextend') for how to cite the package.
##
## Type browseVignettes(package = 'dendextend') for the package vignette.
## The github page is: https://github.com/talgalili/dendextend/
##
## Suggestions and bug-reports can be submitted at: https://github.com/talgalili/dendextend/issues
## Or contact: <tal.galili@gmail.com>
##
## To suppress this message use: suppressPackageStartupMessages(library(dendextend))
## -----
```

```
##
## Attaching package: 'dendextend'
```

```
## The following object is masked from 'package:rpart':
##
##      prune
```

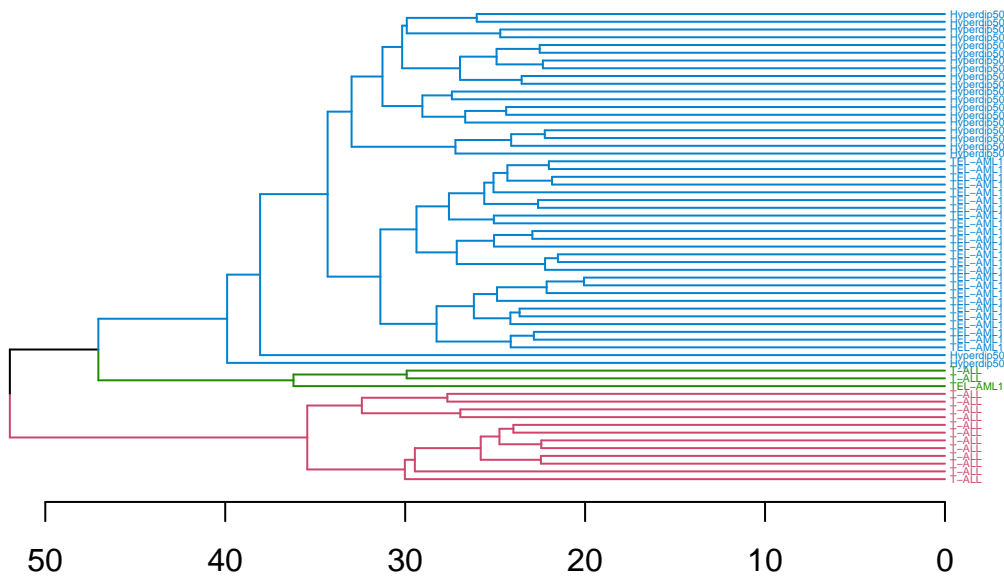
```
## The following object is masked from 'package:stats':
##
##      cutree
```

```
##subsetting to include only rows for which Type is either T-ALL, TEL-AML1, or Hyperdip50
leukemia_subset_1 <- filter(leukemia_data, Type == c("T-ALL","TEL-AML1","Hyperdip50"))
## exclude the first column Type
leukemia_subset = leukemia_subset_1 %>% select(-Type)
## calculate the distance matrix
subset_dist = dist(leukemia_subset)
set.seed(1)
## Hierarchical Clustering using complete linkage
drug.hclust = hclust(subset_dist)
```

```
## first plot
```

```
x = as.dendrogram(drug.hclust)
x %>% set_labels(leukemia_subset_1$Type[order.dendrogram(x)]) %>% set("labels_col",k=3) %>% set("branches_k_color",k=3) %>%
  set("labels_cex", 0.3) %>% plot(main='Three Groups for hclust',horiz = TRUE)
```

### Three Groups for hclust



```
## second plot
```

```
y = as.dendrogram(drug.hclust)
y %>% set_labels(leukemia_subset_1$Type[order.dendrogram(y)]) %>%
  set("labels_col",k=5) %>% set("branches_k_color", k = 5) %>%
  set("labels_cex", 0.3) %>% plot(main='Five Groups for hclust',horiz = TRUE)
```

## Five Groups for hclust

