

Assignment 2

Johnny Lee, s1687781

Problem 1 (27 points)

File wdbc2.csv (available from the accompanying zip folder on Learn) refers to a study of breast cancer where the outcome of interest is the type of the tumour (benign or malignant, recorded in column “diagnosis”). The study collected 30 imaging biomarkers on 569 patients.

Problem 1.a (7 points)

Using package caret, create a data partition so that the training set contains 70% of the observations (set the random seed to 984065 beforehand). Fit both a ridge regression model and a lasso model which uses cross-validation on the training set to diagnose the type of tumour from the 30 biomarkers. Then use a plot to help identify the penalty parameter λ that maximizes the AUC. Note: There is no need to use the prepare.glmnet() function from lab 4, using as.matrix() with the required columns is sufficient.

Answer

```
breast <- read.csv("data_assignment2/wdbc2.csv")

set.seed(984065)

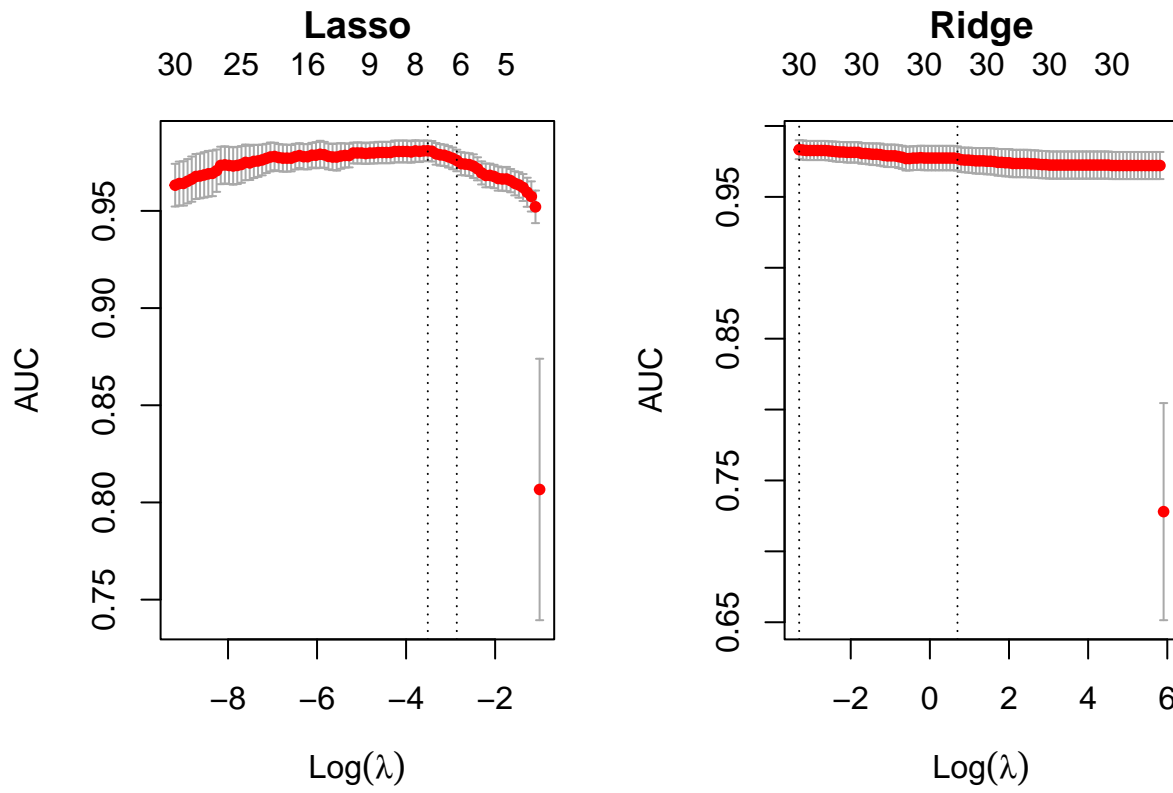
#Splitting into training and testing data set
breast$diagnosis <- ifelse(breast$diagnosis=="benign", 0, 1)
split.index <- createDataPartition(breast$diagnosis,
                                   p = .7, list = TRUE)$Resample1

#train and test data sets
train.breast <- breast[split.index, ]
test.breast <- breast[-split.index, ]

#define explanatory and response variable matrix
biomarkers.x <- as.matrix(subset(train.breast, select = -c(id, diagnosis)))
biomarkers.y <- as.matrix(subset(train.breast, select = c(diagnosis)))

set.seed(984065)
#Fitting ridge regression on training data
fit.lasso <- cv.glmnet(biomarkers.x, biomarkers.y , alpha = 1, family = "binomial", type.measure = "auc")
#Fitting lasso regression on training data
fit.ridge <- cv.glmnet(biomarkers.x, biomarkers.y , alpha = 0, family = "binomial", type.measure = "auc")

par(mfrow=c(1,2), mar=c(4,4,5,2))
plot(fit.lasso, main="Lasso")
plot(fit.ridge, main="Ridge")
```



```
cat("The optimal value of lambda is:", fit.ridge$lambda.min)
```

```
## The optimal value of lambda is: 0.03677378
```

```
cat("The optimal value of lambda is:", fit.lasso$lambda.min)
```

```
## The optimal value of lambda is: 0.02982835
```

Table 1: Lambda values with its model size and AUC

model	Lambda	ModelSize	AUC
Lasso.min	0.0298	8	0.981
Lasso.1se	0.0572	6	0.976
Ridge.min	0.0368	30	0.983
Ridge.1se	2.0100	30	0.977

Problem 1.b (2 points)

Create a data table that for each value of 'lambda.min' and 'lambda.1se' for each model fitted in problem 1.a reports: * the corresponding AUC, * the corresponding model size. Use 3 significant digits for floating point values and comment on these results. Hint: The AUC values are stored in the field called 'cvm'.

Answer

```
#retrieving lambda min from lasso and ridge
lambdamin.lasso <- fit.lasso$lambda.min
lambdamin.ridge <- fit.ridge$lambda.min
#Position of lambda min
index_lambdamin.lasso <- which(lambdamin.lasso == fit.lasso$lambda)
index_lambdamin.ridge <- which(lambdamin.ridge == fit.ridge$lambda)
#retrieving lambda lse from lasso and ridge
lambdalse.lasso <- fit.lasso$lambda.1se
lambdalse.ridge <- fit.ridge$lambda.1se
#Position of lambda lse
index_lambda1se.lasso <- which(lambdalse.lasso == fit.lasso$lambda)
index_lambda1se.ridge <- which(lambdalse.ridge == fit.ridge$lambda)
#AUC values of each lambda values
AUC.lambdamin.lasso <- signif(fit.lasso$cvm[index_lambdamin.lasso],3)
AUC.lambdalse.lasso <- signif(fit.lasso$cvm[index_lambda1se.lasso],3)
AUC.lambdamin.ridge <- signif(fit.ridge$cvm[index_lambdamin.ridge],3)
AUC.lambdalse.ridge <- signif(fit.ridge$cvm[index_lambda1se.ridge],3)
#Model size of each lambda values
ms.lasso.min <- fit.lasso$nzero[index_lambdamin.lasso]
ms.lasso.1se <- fit.lasso$nzero[index_lambda1se.lasso]
ms.ridge.min <- fit.ridge$nzero[index_lambdamin.ridge]
ms.ridge.1se <- fit.ridge$nzero[index_lambda1se.ridge]
table <- data.table(model = c("Lasso.min", "Lasso.1se",
                             "Ridge.min", "Ridge.1se"),
                   Lambda = c(signif(lambdamin.lasso,3),
                               signif(lambdalse.lasso,3),
                               signif(lambdamin.ridge,3),
                               signif(lambdalse.ridge,3)),
                   ModelSize = c(ms.lasso.min, ms.lasso.1se,
                                 ms.ridge.min, ms.ridge.1se),
                   AUC = c(AUC.lambdamin.lasso, AUC.lambdalse.lasso,
                           AUC.lambdamin.ridge, AUC.lambdalse.ridge))
kable(table, caption = "Lambda values with its model size and AUC")
```

Problem 1.c (7 points)

Perform both backward (we'll later refer to this as model B) and forward (model S) stepwise selection on the same training set derived in problem 1.a. Report the variables selected and their standardized regression coefficients in decreasing order of the absolute value of their standardized regression coefficient. Discuss the results and how the different variables entering or leaving the model influenced the final result.

Answer

The forward model takes a null model(which is a model without any feature) and only considers the intercept as a starting point and then progress towards the full model (with 30 features) by adding features. It can be seen here that the final model (forward) considers more features than the backward model. The forward model considers 15 features instead of 14 (in the backwards model)

```
#Computing coefficients of each model and tabulating it.
B.coef <-model.B$coefficients
B.order<- order(abs(B.coef), decreasing = TRUE)
S.coef <-model.S$coefficients
S.order<- order(abs(S.coef), decreasing = TRUE)
table <- list(B.coef[B.order], S.coef[S.order])
kable(table, col.names = "Coefficients", caption = "Coefficients of Model B and Model S") %>%
  kable_styling(latex_options = "hold_position")
```

Table 2: Coefficients of Model B and Model S

	Coefficients		Coefficients
smoothness.stderr	16.9805981	concavepoints	5.8623031
concavepoints	5.3127566	fractaldimension.worst	3.3772683
fractaldimension.worst	4.7503948	compactness	-2.9706608
compactness.stderr	-3.8819567	smoothness.worst	2.4703766
compactness	-2.5826629	(Intercept)	-2.2713307
(Intercept)	-2.3003927	radius.stderr	0.9152699
concavity.stderr	0.9488969	symmetry.worst	0.7768531
symmetry.worst	0.8536208	concavepoints.worst	-0.7037566
radius.stderr	0.4256226	concavity.worst	0.4276025
concavity.worst	0.3504861	radius.worst	0.1106887
radius.worst	0.1408605	radius	0.0856111
radius	0.0743984	perimeter.stderr	-0.0503778
perimeter	-0.0111032	texture.worst	0.0108587
texture.worst	0.0105318	perimeter	-0.0102155
area.worst	-0.0009868	area.worst	-0.0008448
		area.stderr	-0.0007198

Problem 1.d (3 points)

Compare the goodness of fit of model B and model S in an appropriate way.

Answer

```
#Testing goodness of fit of model B and model S  
cat("Model B deviance: ", model.B$deviance)
```

```
## Model B deviance: 24.59436
```

```
cat("Model S deviance: ", model.S$deviance)
```

```
## Model S deviance: 24.57668
```

```
pchisq(model.B$null.deviance - model.B$deviance, df = 15, lower.tail = FALSE)
```

```
## [1] 7.40918e-09
```

```
pchisq(model.S$null.deviance - model.S$deviance, df = 15, lower.tail = FALSE)
```

```
## [1] 7.355861e-09
```

Problem 1.e (2 points)

Compute the training AUC for model B and model S.

Answer

```
model.B.auc <- roc(train.breast$diagnosis, model.B$fitted.values, plot = TRUE,  
  xlim = c(0,1), col="red")
```

```
## Setting levels: control = 0, case = 1
```

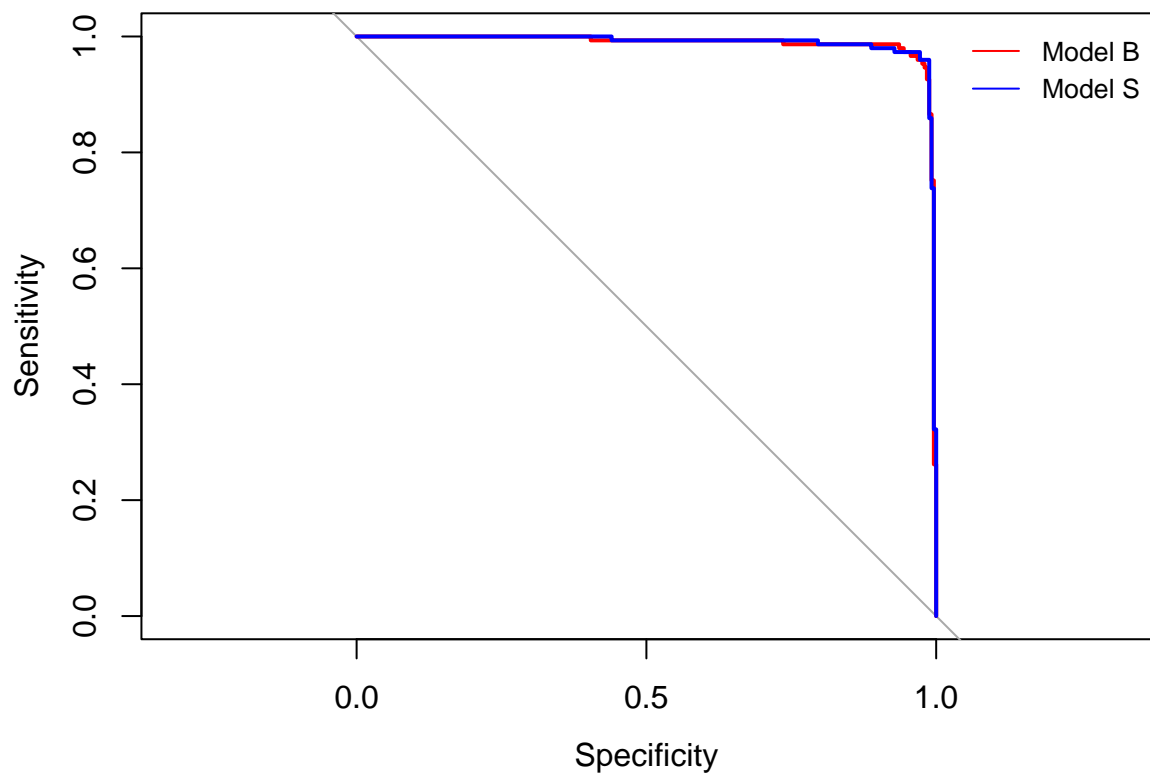
```
## Setting direction: controls < cases
```

```
model.S.auc <- roc(train.breast$diagnosis, model.S$fitted.values, plot = TRUE,  
  add = TRUE,col = "blue")
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
legend("topright", legend = c("Model B", "Model S"),  
  col = c("red", "blue"), lty = 1, cex = 0.8, bty = "n")
```



Problem 1.f (6 points)

Use the four models to predict the outcome for the observations in the test set (use the lambda at 1 standard error for the penalised models). Plot the ROC curves of these models (on the sameplot, using different colours) and report their test AUCs. Compare the training AUCs obtained in problems 1.b and 1.e with the test AUCs and discuss the fit of the different models.

Answer

```
#Lasso model
lasso.pred <- predict(fit.lasso, newx = as.matrix(test.breast[, -c(1,2)]), s="lambda.1se")
#Ridge Regression model
ridge.pred <- predict(fit.ridge, newx = as.matrix(test.breast[, -c(1,2)]), s="lambda.1se")
#Model B
B.pred <- predict(model.B, newdata = test.breast, type = "response")
#Model S
S.pred <- predict(model.S, newdata = test.breast, type = "response")
#Plotting ROC
lasso.auc <- roc(test.breast$diagnosis, lasso.pred, plot = TRUE, xlim = c(0,1), col = "red")$auc

## Setting levels: control = 0, case = 1
## Warning in roc.default(test.breast$diagnosis, lasso.pred, plot = TRUE, xlim =
## c(0, : Deprecated use a matrix as predictor. Unexpected results may be produced,
## please pass a numeric vector.
## Setting direction: controls < cases
ridge.auc <- roc(test.breast$diagnosis, ridge.pred, plot = TRUE, col = "blue", add = TRUE)$auc

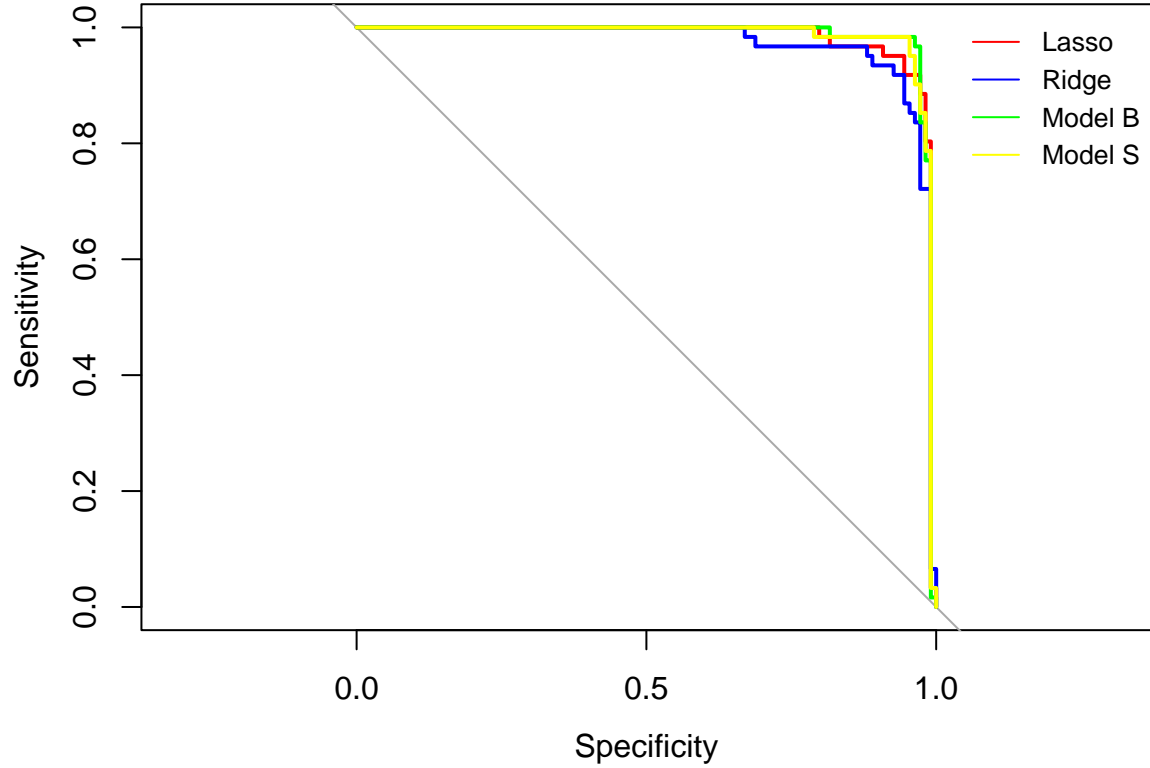
## Setting levels: control = 0, case = 1
## Warning in roc.default(test.breast$diagnosis, ridge.pred, plot = TRUE, col
## = "blue", : Deprecated use a matrix as predictor. Unexpected results may be
## produced, please pass a numeric vector.
## Setting direction: controls < cases
B.auc <- roc(test.breast$diagnosis, B.pred, plot = TRUE, col = "green", add = TRUE)$auc

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
S.auc <- roc(test.breast$diagnosis, S.pred, plot = TRUE, col = "yellow", add = TRUE)$auc

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
legend("topright", legend = c("Lasso", "Ridge", "Model B", "Model S"),
      col = c("red", "blue", "green", "yellow"), lty = 1, cex = 0.8, bty = "n")
```

Table 3: Training and Testing Acuraccy of each Model

models	train.auc	test.auc
Lasso Regression	0.9760000	0.9808994
Ridge Regression	0.9770000	0.9712739
Model B	0.9886980	0.9846593
Model S	0.9894228	0.9837570



```
#Compare the AUCs
train.auc <- c(AUC.lambda1se.lasso, AUC.lambda1se.ridge, model.B.auc$auc, model.S.auc$auc)
test.auc <- c(lasso.auc, ridge.auc, B.auc, S.auc)
models <- c("Lasso Regression", "Ridge Regression", "Model B", "Model S")
# We make a table and voila!
table <- data.table(models, train.auc, test.auc)
kable(table, caption = "Training and Testing Acuraccy of each Model")
```


Problem 2 (40 points)

File GDM.raw.txt (available from the accompanying zip folder on Learn) contains 176 SNPs to be studied for association with incidence of gestational diabetes (a form of diabetes that is specific to pregnant women). SNP names are given in the form “rs1234_X” where “rs1234” is the official identifier (rsID), and “X” (one of A, C, G, T) is the reference allele.

Problem 2.a (3 points)

Read file GDM.raw.txt into a data table named gdm.dt. Impute missing values in gdm.dt according to SNP-wise median allele count.

Answer

```
gdm.dt <- data.table(fread("data_assignment2/GDM.raw.txt"))

#Performing Imputation using median
for (colnm in colnames(gdm.dt)[-1]){
  gdm.dt[[colnm]][is.na(gdm.dt[[colnm]])] <- median(gdm.dt[[colnm]], na.rm = T)
}
```

Problem 2.b (8 points)

Write function `univ.glm.test <- function(x, y, order = FALSE)` where `x` is a data table of SNPs, `y` is a binary outcome vector, and `order` is a boolean. The function should fit a logistic regression model for each SNP in `x`, and return a data table containing SNP names, regression coefficients, odds ratios, standard errors and p-values. If `order` is set to `TRUE`, the output data table should be ordered by increasing p-value.

Answer

```
univ.glm.test <- function(x, y, order = FALSE){
  stopifnot(nrow(x) == length(y))
  output <- data.table("SNP"=character(),
                       "coefficients"=numeric(), "odds.ratios"=numeric(),
                       "std.error"=numeric(), "p.value"=numeric())
  for (i in 1:ncol(x)){
    regr <- glm(y ~ x[[i]], family = binomial(link = "logit"))
    summarised <- coef(summary(regr))
    ## remove the row corresponding to the intercept and the column containing
    ## the t-value, then convert to a dataframe
    output <- rbind(output, list(names(x)[i], summarised[2,1],
                                exp(summarised[2,1]), summarised[2,2],
                                summarised[2,4]))
  }
  ## assign better column names
  # colnames(output) <- c("SNP", "Coefficients", "odds ratios", "std.error", "p.value")
  if(order == TRUE){
    output <- output[order(p.value)]
  }
  return(output)
}
```

Problem 2.c (5 points)

Using function `univ.glm.test()`, run an association study for all the SNPs in `gdm.dt` against having gestational diabetes (column “pheno”). For the SNP that is most strongly associated to increased risk of gestational diabetes and the one with most significant protective effect, report the summary statistics from the GWAS as well as the 95% and 99% confidence intervals on the odds ratio.

Answer

```
x <- gdm.dt[, 4:ncol(gdm.dt)]
y <- gdm.dt[[3]]
study <- univ.glm.test(x, y)
kable(head(study), caption = "Logistic Regression on Pheno vs each SNP") %>%
  kable_styling(latex_options = "hold_position")
```

Table 4: Logistic Regression on Pheno vs each SNP

SNP	coefficients	odds.ratios	std.error	p.value
rs7513574_T	0.0021575	1.002160	0.1051372	0.9836280
rs1627238_A	0.1146379	1.121467	0.1138224	0.3138559
rs1171278_C	0.1214094	1.129087	0.1138073	0.2860628
rs1137100_A	0.0601048	1.061948	0.1104238	0.5862285
rs2568958_A	0.1493799	1.161114	0.1233800	0.2259989
rs1514175_A	0.0562296	1.057841	0.1052359	0.5931203

```
index <- which(study$odds.ratio == max(study$odds.ratio))
kable(study[index, ],
      caption=" most strongly associated to increased risk of gestational diabetes") %>%
  kable_styling(full_width = F, position = "center", latex_options = "hold_position")
```

Table 5: most strongly associated to increased risk of gestational diabetes

SNP	coefficients	odds.ratios	std.error	p.value
rs1423096_T	0.6510641	1.91758	0.3166547	0.0397758

```
strong.coef <- study[index, ]$coefficients
strong.se <- study[index, ]$std.error

confidence.int.95 <- round(exp(strong.coef + 1.96 * strong.se*c(-1,1)), 3)
confidence.int.99 <- round(exp(strong.coef + 2.576 * strong.se*c(-1,1)),3)

newindex <- which(study$odds.ratio < 1)
best <- study[newindex,]
# Select the SNP with lowest p value
index3 <- which(best$p.value == min(best$p.value))
best.SNP <- best[index3]
kable(best.SNP,
      caption="most protective effect on gestational diabetes") %>%
  kable_styling(full_width = F, position = "center", latex_options = "hold_position")

#Compute Confidence Interval
best.coef <- best.SNP$coefficients
best.se <- best.SNP$std.error
```

Table 6: most protective effect on gestational diabetes

SNP	coefficients	odds.ratios	std.error	p.value
rs2237897_T	-0.4394456	0.6443936	0.1126133	9.53e-05

```
best.ci.95 <- round(exp(best.coef +1.96*best.se *c(-1,1)),3)
best.ci.99 <- round(exp(best.coef +2.576*best.se *c(-1,1)),3)

#Output
cat(" SNP most strongly associated to increased risk of gestational      diabetes, ",
    "\n 95% Confidence Interval is", confidence.int.95,
    "\n 99% Confidence Interval is", confidence.int.99)

## SNP most strongly associated to increased risk of gestational      diabetes,
## 95% Confidence Interval is 1.031 3.567
## 99% Confidence Interval is 0.848 4.335

cat("\n SNPs with most protective effect on gestational diabetes,",
    "\n 95% Confidence Interval is", best.ci.95,
    "\n 99% Confidence Interval is", best.ci.99)

##
## SNPs with most protective effect on gestational diabetes,
## 95% Confidence Interval is 0.517 0.804
## 99% Confidence Interval is 0.482 0.861
```

We can see that SNP rs1423096_T has the highest odds ratio (1.91758) and hence is the most strongly associated to increased risk of gestational diabetes. In fact, this SNP increases the odds of having gestational diabetes by about 92%. The SNP with most significant protective effect is rs2237897_T and it reduced the risk of diabetes by about 35%.

Problem 2.d (4points)

Merge your GWAS results with the table of gene names provided in file GDM.annot.txt (available from the accompanying zip folder on Learn). For SNPs that have p-value $< 10^{-4}$ (hit SNPs) report SNP name, effect allele, chromosome number and corresponding gene name. Separately, report for each 'hit SNP' the names of the genes that are within a 1Mb window from the SNP position on the chromosome. Note: That's genes that fall within $\pm 1,000,000$ positions using the 'pos' column in the dataset.

Answer

```
gene.names <- fread("data_assignment2/GDM.annot.txt")
# study
# new.study <- study
new.study <- study %>% copy() %>%
  .[, full.snp := SNP] %>%
  .[, c("SNP", "effect.allele") := do.call(Map, c(f = c, strsplit(SNP, "_")))]

merge.dt <- merge(new.study, gene.names, by.x = "SNP", by.y="snp", all = TRUE)
hit.snps <- merge.dt[p.value<1e-4]
kable(hit.snps[,c("SNP", "effect.allele","chrom","gene")],
      caption= "SNPs that have p-value  $< 10^{-4}$ ") %>%
  kable_styling(latex_options = "hold_position")
```

Table 7: SNPs that have p-value $< 10^{-4}$

SNP	effect.allele	chrom	gene
rs12243326	A	10	TCF7L2
rs2237897	T	11	KCNQ1

```

hit.snp.window <- data.table()
for (i in hit.snps$SNP){
  idx <- which(hit.snps$SNP == i)
  window.val <- merge.dt[(merge.dt$pos>= hit.snps$pos[idx] - 1e6) &
                        (merge.dt$pos<= hit.snps$pos[idx] + 1e6)]
  hit.snp.window <- rbind(hit.snp.window, window.val)
}
# Display the genes that fall within this window
kable(data.table(hit.snp.window$gene), col.names = "gene",
      caption = "Gene within 1Mb Window") %>%
kable_styling(latex_options = "hold_position")

```

Table 8: Gene within 1Mb Window

gene
TCF7L2
TCF7L2
TCF7L2
TCF7L2
TH
KCNQ1
CACNA2D4
KCNQ1
KCNQ1
KCNQ1
SMG6
SMG6

Problem 2.e (8 points)

Build a weighted genetic risk score that includes all SNPs with p-value $< 10^{-4}$, a score with all SNPs with p-value $< 10^{-3}$, and a score that only includes SNPs on the FTO gene (hint: ensure that the ordering of SNPs is respected). Add the three scores as columns to the gdm.dt data table. Fit the three scores in separate logistic regression models to test their association with gestational diabetes, and for each report odds ratio, 95% confidence interval and p-value.

Answer

```

#Defining each weighted genetic risk score
hit.snp.1 <- merge.dt[p.value<1e-4]
gdm.1 <- gdm.dt[, .SD, .SDcols = merge.dt[p.value <1e-4]$full.snp]
names(gdm.1) <- gsub("_.", "", x=names(gdm.1))
wgrs.1 <- as.matrix(gdm.1) %*% hit.snp.1$coefficients

hit.snp.2 <- merge.dt[p.value<1e-3]

```

```

gdm.2 <- gdm.dt[, .SD, .SDcols = merge.dt[p.value <1e-3]$full.snp]
names(gdm.2) <- gsub("_.", "", x=names(gdm.2))
wgrs.2 <- as.matrix(gdm.2) %*% hit.snp.2$coefficients

hit.snp.3 <- merge.dt[gene=="FTO"]
gdm.3 <- gdm.dt[, .SD, .SDcols = merge.dt[gene=="FTO"]$full.snp]
names(gdm.3) <- gsub("_.", "", x=names(gdm.3))
wgrs.3 <- as.matrix(gdm.3) %*% hit.snp.3$coefficients

#Adding 3 columns to gdm.dt
scores <- c("score.1", "score.2", "score.3")
gdm.dt <- gdm.dt %>% copy() %>%
  .[, `:=`(scores.1=wgrs.1, scores.2=wgrs.2, scores.3=wgrs.3)]

y <- gdm.dt[[3]]
x <- gdm.dt[,180:182]
wgrs.snp <- univ.glm.test(x, y)
wgrs.snp <- wgrs.snp %>%
  .[,upper.conf.int:=round(exp(coefficients + 1.96 * std.error*-1), 3)] %>%
  .[,lower.conf.int:=round(exp(coefficients + 1.96 * std.error), 3)] %>%
  .[, !"coefficients"] %>% .[, !"std.error"]

kable(head(wgrs.snp), caption = "Logistic regression on Pheno vs SNP") %>%
  kable_styling(latex_options = "hold_position")

```

Table 9: Logistic regression on Pheno vs SNP

SNP	odds.ratios	p.value	upper.conf.int	lower.conf.int
scores.1	2.729433	0.0000000	1.915	3.890
scores.2	1.451854	0.0000000	1.279	1.648
scores.3	1.413857	0.2151883	0.818	2.445

Problem 2.f (4 points)

File GDM.test.txt (available from the accompanying zip folder on Learn) contains genotypes of another 40 pregnant women with and without gestational diabetes (assume that the reference allele is the same one that was specified in file GDM.raw.txt). Read the file into variable gdm.test. For the set of patients in gdm.test, compute the three genetic risk scores as defined in problem 2.e using the same set of SNPs and corresponding weights. Add the three scores as columns to gdm.test (hint: use the same columnnames as before).

Answer

```

gdm.test <- setDT(fread("data_assignment2/GDM.test.txt"))

gdm1.colname <- colnames(gdm.1)
gdm.test.1 <- gdm.test[,..gdm1.colname]
wgrs.test.1 <- as.matrix(gdm.test.1) %*% hit.snp.1$coefficients
gdm2.colname <- colnames(gdm.2)
gdm.test.2 <- gdm.test[,..gdm2.colname]
wgrs.test.2 <- as.matrix(gdm.test.2) %*% hit.snp.2$coefficients
gdm3.colname <- colnames(gdm.3)
gdm.test.3 <- gdm.test[,..gdm3.colname]
wgrs.test.3 <- as.matrix(gdm.test.3) %*% hit.snp.3$coefficients

```

```
#Adding 3 columns to gdm.test
scores <- c("score.1", "score.2", "score.3")
gdm.test <- gdm.test %>% copy() %>%
  .[,`:=`(scores.1 = wgrs.test.1, scores.2 = wgrs.test.2, scores.3 = wgrs.test.3)]
```

Problem 2.g (4 points)

Use the logistic regression models fitted in problem 2.e to predict the outcome of patients in gdm.test. Compute the test log-likelihood for the predicted probabilities from the three genetic risk score models.

Answer

```
fit1 <- glm(y ~ scores.1, family = binomial(link = "logit"), data=gdm.dt)
pred.1 <- predict(fit1, newdata = gdm.test[,180], type="response")
fit2 <- glm(y ~ scores.2, family = binomial(link = "logit"), data=gdm.dt)
pred.2 <- predict(fit2, newdata = gdm.test[,181], type="response")
fit3 <- glm(y ~ scores.3, family = binomial(link = "logit"), data=gdm.dt)
pred.3 <- predict(fit3, newdata = gdm.test[,182], type="response")
cat("The log likelihood of score 1:",
    sum(gdm.test$pheno*log(pred.1)+(1-gdm.test$pheno)*log(1-pred.1)))
```

```
## The log likelihood of score 1: -25.06824
```

```
cat("The log likelihood of score 1:",
    sum(gdm.test$pheno*log(pred.2)+(1-gdm.test$pheno)*log(1-pred.2)))
```

```
## The log likelihood of score 1: -24.77693
```

```
cat("The log likelihood of score 1:",
    sum(gdm.test$pheno*log(pred.3)+(1-gdm.test$pheno)*log(1-pred.3)))
```

```
## The log likelihood of score 1: -28.05355
```

Problem 2.h (4points)

File GDM.study2.txt (available from the accompanying zip folder on Learn) contains the summary statistics from a different study on the same set of SNPs. Perform a meta-analysis with the results obtained in problem 2.c (hint: remember that the effect alleles should correspond) and produce a summary of the meta-analysis results for the set of SNPs with meta-analysis p-value $< 10^{-4}$ sorted by increasing p-value.

Answer

```
gdm.study <- setDT(fread("data_assignment2/GDM.study2.txt"))

are.equal <- gdm.study$effect.allele == new.study$effect.
not.equal <- gdm.study$effect.allele != new.study$effect.
kable(table(are.equal, not.equal), caption = "Confusion Matrix") %>%
  kable_styling(latex_options = "hold_position")
```

Table 10: Confusion Matrix

	FALSE	TRUE
FALSE	0	141
TRUE	35	0

```

beta1 <- gdm.study$beta
beta2 <- new.study$coefficients
beta2[not.equal] <- -beta2[not.equal]

weight1 <- 1/ gdm.study$se
weight2 <- 1/ new.study$std.error
head(weight1)

## [1] 5.901481 2.891623 1.182123 6.207614 3.358099 3.023486

head(weight2)

## [1] 9.511380 8.785615 8.786784 9.056015 8.105040 9.502461

beta_meta_analysis = (weight1*beta1 + weight2*beta2)/(weight1 + weight2)
standard_error_meta_analysis = sqrt(1 / weight1 + weight2)
p_value_meta_analysis = 2*pnorm(abs(beta_meta_analysis / standard_error_meta_analysis), lower.tail = F)
summary = merge(gdm.study, new.study, by.x = c("snp","effect.allele"),by.y=c("SNP","effect.allele"), all=TRUE)
# [,c("snp", "effect.allele", "other.allele")]
summary = cbind(summary, data.table(beta_meta_analysis = beta_meta_analysis, standard_error_meta_analysis = standard_error_meta_analysis))

## Warning in as.data.table.list(x, keep.rownames = keep.rownames, check.names
## = check.names, : Item 2 has 176 rows but longest item has 205; recycled with
## remainder.

```

Problem 3 (33 points)

File nki.csv (available from the accompanying zip folder on Learn) contains data for 144 breast cancer patients. The dataset contains a binary outcome variable (“Event”, indicating the insurgence of further complications after operation), covariates describing the tumour and the age of the patient, and gene expressions for 70 genes found to be prognostic of survival.

Problem 3.a (6 points)

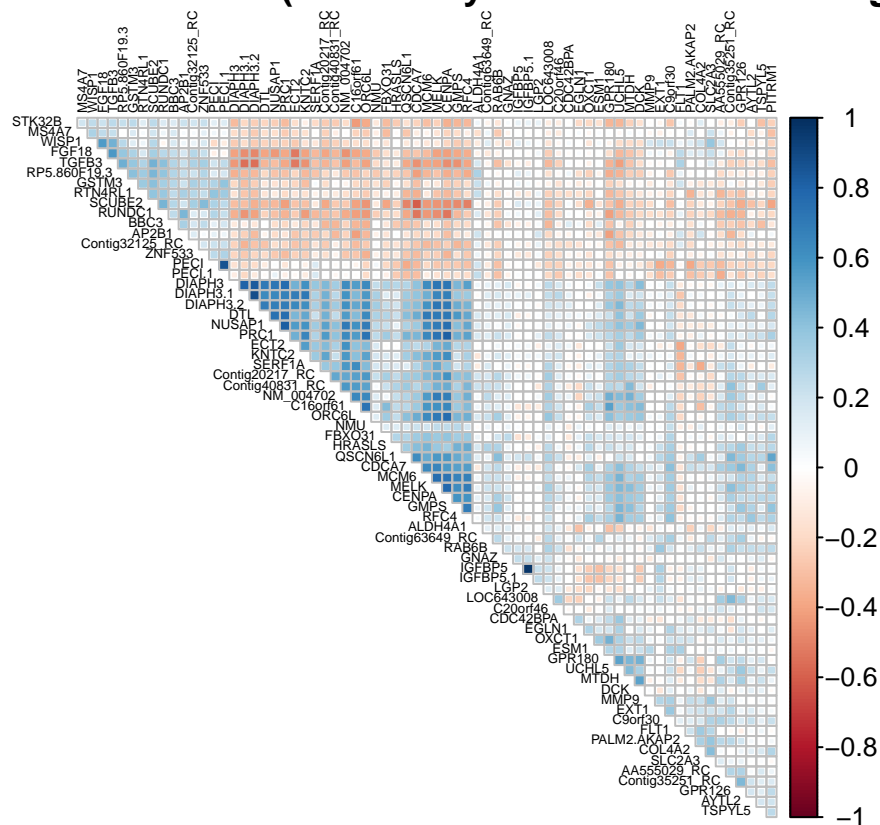
Compute the matrix of correlations between the gene expression variables, and display it so that a block structure is highlighted. Discuss what you observe. Write some code to identify the unique pairs of (distinct) variables that have correlation coefficient greater than 0.80 in absolute value and report their correlation coefficients.

Answer

```
nki <- read.csv("data_assignment2/nki.csv")

nki.cor <- cor(nki[,7:76], use="pairwise.complete")
corrplot(nki.cor, order="hclust", diag=FALSE, tl.col="black", tl.cex = 0.4, method = "square", title="")
```

Correlation matrix (ordered by hierarchical clustering)



```
gene1 <- gene2 <- corr <- c()
for (i in 1:nrow(nki.cor)){
  for (j in 1:ncol(nki.cor)){
    if (abs(nki.cor[i,j])>0.8 & nki.cor[i,j] !=1){
      gene1 <- c(gene1, rownames(nki.cor)[i])
      gene2 <- c(gene2, colnames(nki.cor)[j])
    }
  }
}
```



```

      corr <- c(corr, nki.cor[i,j])
    }
  }
}

cor0.8 <- data.table(gene1, gene2, corr)
kable(unique(cor0.8, by = "corr"), caption = "Correlation Coefficient  $> 0.8$ ") %>%
  kable_styling(latex_options = "hold_position")

```

Table 11: Correlation Coefficient > 0.8

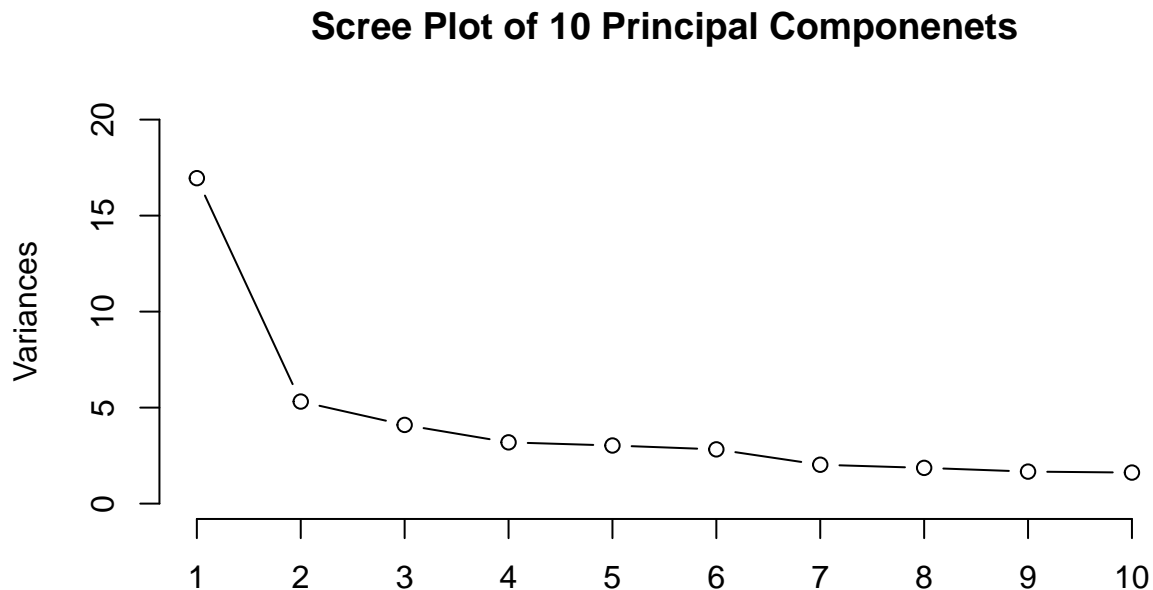
gene1	gene2	corr
DIAPH3	DIAPH3.1	0.8031368
DIAPH3	DIAPH3.2	0.8338591
NUSAP1	PRC1	0.8298356
DIAPH3.1	DIAPH3.2	0.8868741
PECI	PECI.1	0.8697836
IGFBP5	IGFBP5.1	0.9775030
PRC1	CENPA	0.8175424

Problem 3.b (8 points)

Run PCA (only over the columns containing gene expressions), in order to derive a patient-wise summary of all gene expressions (dimensionality reduction). Decide which components to keep and justify your decision. Test if those principal components are associated with the outcome in unadjusted logistic regression models and in models adjusted for age, estrogen receptor and grade. Justify the difference in results between unadjusted and adjusted models.

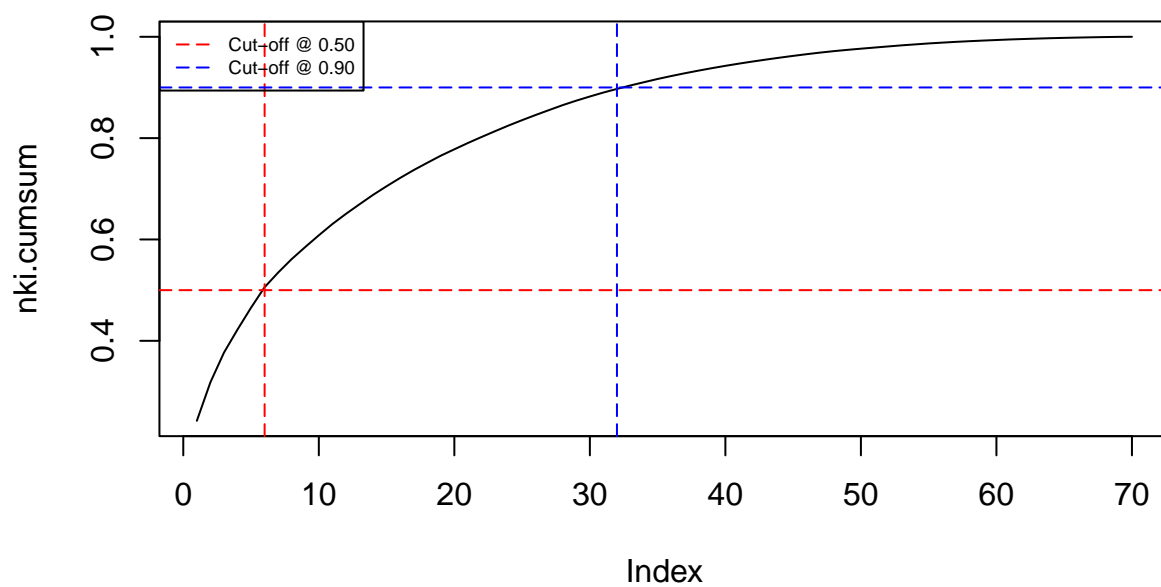
Answer

```
nki.pca<- prcomp(nki[,7:76], center=TRUE, scale = TRUE)
nki.cumsum <- cumsum(nki.pca$sdev^2 / sum(nki.pca$sdev^2))
plot(nki.pca, type="line", main="Scree Plot of 10 Principal Componenets",
     ylim = c(0,20))
```

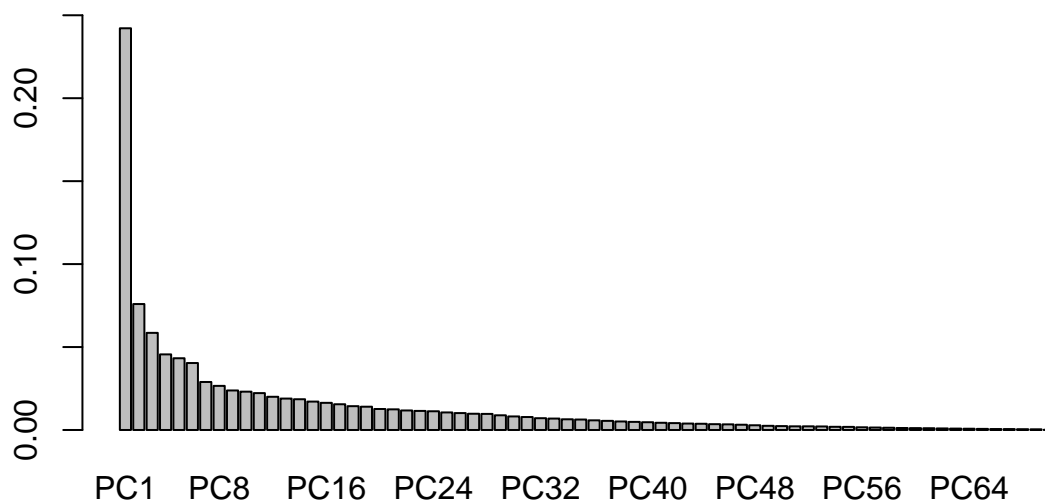


```
plot(nki.cumsum, type="line", main="Cumulative Variance of Principal Components")
abline(v = 6, col="red", lty=5)
abline(v = 32, col="blue", lty=5)
abline(h = 0.50, col="red", lty=5)
abline(h = 0.90, col="blue", lty=5)
legend("topleft", legend=c("Cut-off @ 0.50", "Cut-off @ 0.90"),
     col=c("red","blue"), lty=5, cex=0.6)
```

Cumulative Variance of Principal Components



```
barplot(summary(nki.pca)$importance[2,], ylim=c(0, 0.25))
```

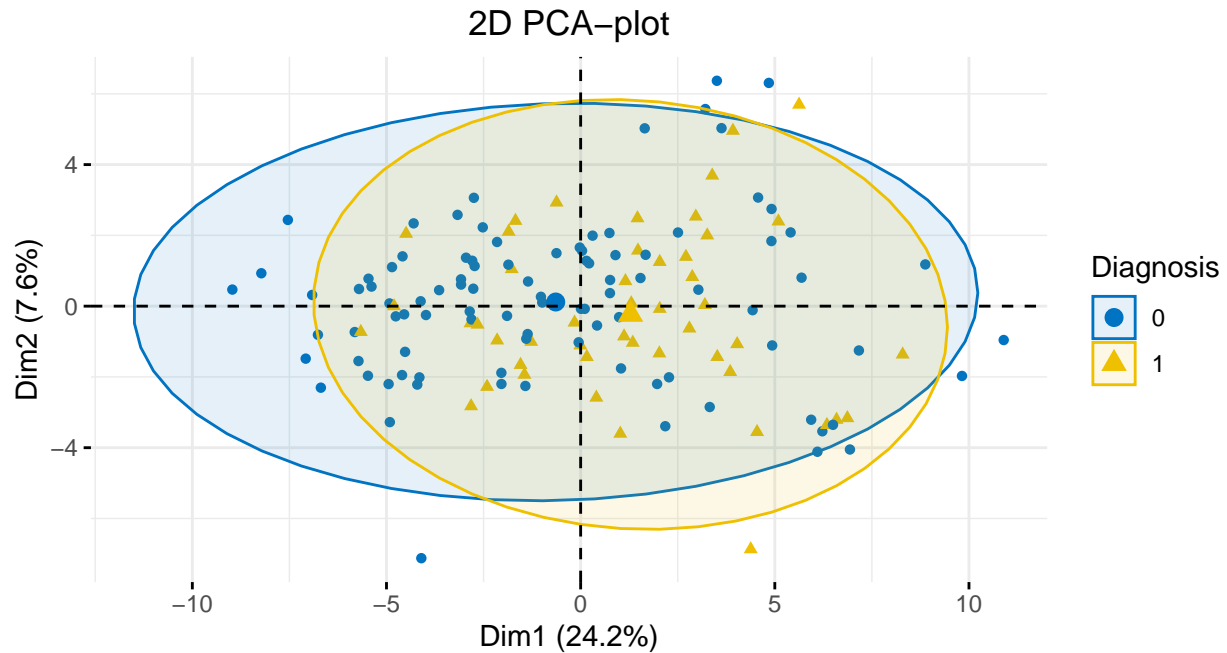


Problem 3.c (8 points)

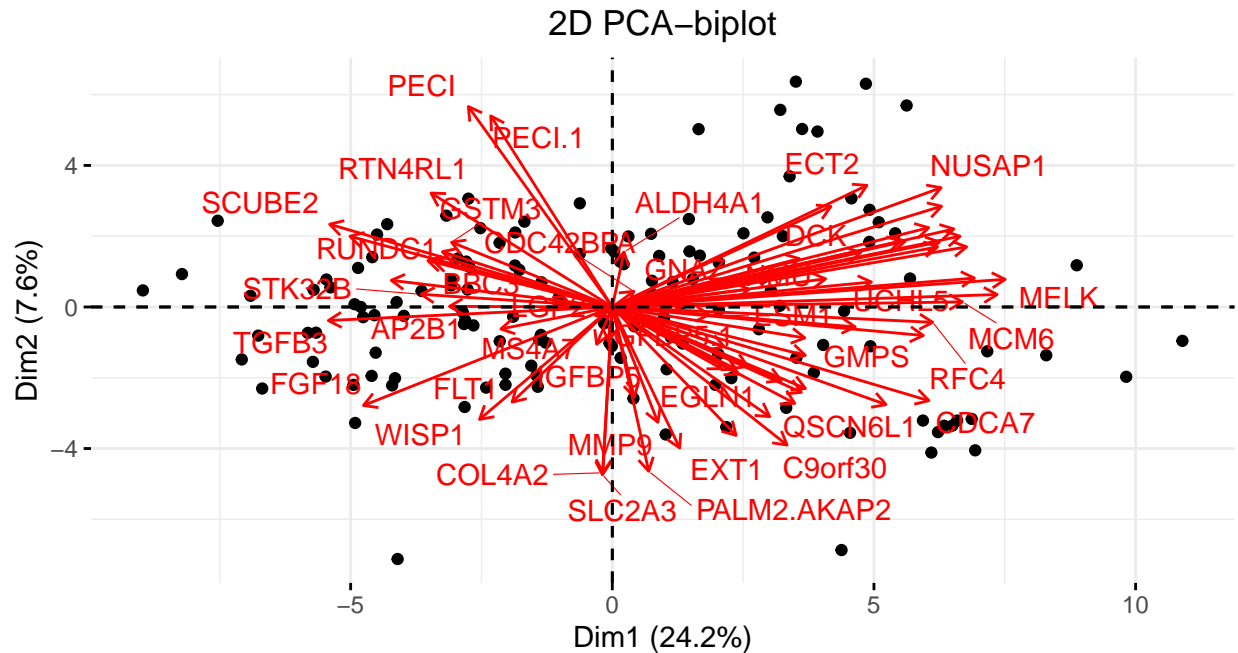
Use plots to compare with the correlation structure observed in problem 2.a and to examine how well the dataset may explain your outcome. Discuss your findings and suggest any further steps if needed.

Answer

```
fviz_pca_ind(nki.pca, geom='point', habillage = nki$Event, axes = c(1,2),
             addEllipses = T, palette="jco", legend.title="Diagnosis") +
  ggtitle("2D PCA-plot") +
  theme(plot.title = element_text(hjust = 0.5))
```



```
fviz_pca_biplot(nki.pca, geom = 'point', repel = T, col.var = "red") +
  ggtitle("2D PCA-biplot") +
  theme(plot.title = element_text(hjust = 0.5))
```



Problem 3.d (11 points)

Based on the models we examined in the labs, fit an appropriate model with the aim to provide the most accurate prognosis you can for patients. Discuss and justify your decisions.

Answer

```
invisible({capture.output({
  set.seed(984065)
  event <- nki$Event
  split.index <- createDataPartition(event, p = 0.7)$Resample1
  gene.response <- as.matrix(nki[, 1:6])

  #Training and testing data set
  nki.x.train <- gene.response[split.index,]
  nki.x.test <- gene.response[-split.index,]
  nki.y.train <- event[split.index]
  nki.y.test <- event[-split.index]

  nki.train = nki[split.index,]
  nki.test = nki[-split.index, ]

  nki.full <- glm(Event ~., data = nki.train, family = binomial(link = "logit"))
  nki.null <- glm(Event ~ 1, data = nki.train, family = binomial(link = "logit"))
  #Forward and backward models
  nki.forward <- stepAIC(nki.null, scope=list(upper=nki.full),
                        direction = "forward")
  nki.backward <- stepAIC(nki.full, scope=list(upper=nki.null),
                        direction = "backward")
}))})
```

```
nki.forward$aic
```

```
## [1] 44
```

```
as.formula(nki.forward)
```

```
## Event ~ LymphNodes + NUSAP1 + UCHL5 + Contig63649_RC + QSCN6L1 +  
##      Contig20217_RC + HRASLS + Contig32125_RC + STK32B + LGP2 +  
##      PRC1 + LOC643008 + EGLN1 + CDCA7 + ECT2 + KNTC2 + Diam +  
##      GPR180 + RTN4RL1 + IGFBP5 + PITRM1
```

```
nki.backward$aic
```

```
## [1] 44
```

```
as.formula(nki.backward)
```

```
## Event ~ Diam + LymphNodes + Contig63649_RC + FGF18 + DIAPH3.1 +  
##      Contig32125_RC + C16orf61 + OXCT1 + MMP9 + KNTC2 + UCHL5 +  
##      STK32B + DCK + SLC2A3 + PECI.1 + LOC643008 + MCM6 + IGFBP5 +  
##      HRASLS + LGP2 + EGLN1
```