

Assignment 2

Biomedical Data Science (MATH11174), 22/23, Semester 2

April 6, 2023

Due on Thursday, 6th of April 2023, 5:00pm

! Pay Attention

The assignment is marked out of 100 points, and will contribute to *30%* of your final mark. The aim of this assignment is to produce a precise report in biomedical studies with the help of statistical and machine learning. Please complete this assignment using **Quarto/Rmarkdown file and render/knit this document only in PDF format** (rendering while solving the questions will prevent sudden panic before submission!). Submit using the **gradescope link on Learn** and ensure that **all questions are tagged accordingly**. You can simply click render on the top left of Rstudio (**Ctrl+Shift+K**). If you cannot render/knit to PDF directly, open Terminal in your RStudio (**Alt+Shift+R**) and type **quarto tools install tinytex**, otherwise please follow this [link](#). If you have any code that does not run you will not be able to render nor knit the document so comment it as you might still get some grades for partial code.

Codes that are **clear and reusable** will be rewarded. Codes without proper indentation, choice of variable identifiers, **comments**, efficient code, etc will be penalised. An initial code chunk is provided after each subquestion but **create as many chunks as you feel is necessary** to make a clear report. Add plain text explanations in between the chunks when required to make it easier to follow your code and reasoning. Ensure that all answers containing multiple values should be presented and formatted only with **kable()** and **kable_styling()** otherwise penalised (no use of **print()** or **cat()**). All plots must be displayed with clear title, label and legend otherwise penalised.

This is an **individual assignment**, and **no public discussions** will be allowed. If you have any question, please ask on Piazza by specifying your **Post to** option to **instructors**. To join Piazza, please follow this [link](#).

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 and Lasso regression model which use 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 maximises the AUC and report the λ for both ridge and Lasso regression using `kable()`.
- *Note : there is no need to use the `prepare.glmnet()` function from lab 4, using `as.matrix()` with the required columns is sufficient.*

```

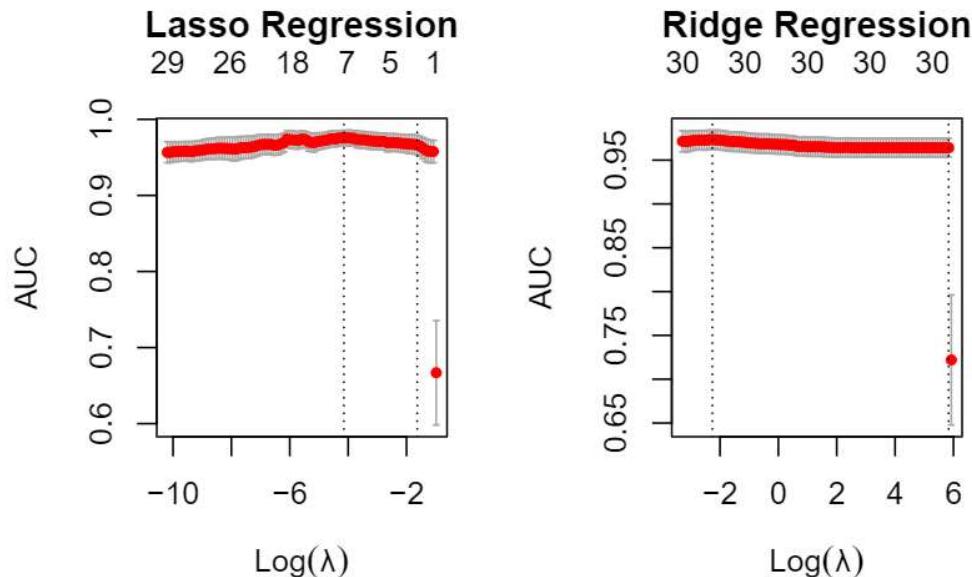
1 # s2319494
2 # Set given seed value using the set.seed() function.
3 set.seed(984065)
4 # Read the dataset wdbc2.csv using fread().
5 wdbc.dt<-fread('wdbc2.csv')
6 # Convert given dataset to Data table
7 wdbc.dt<-setDT(wdbc.dt)
8 # Now it is time to do a train-test split of 70:30. This can be
9 # done using the createDataPartition() function specifying p=0.7.
10 # The data is split once and thus times=1 and we don't have to
11 # view the result as a list and hence list=FALSE.
12 t_seti<-createDataPartition(wdbc.dt$diagnosis,
13 #           p=0.7,
14 #           list = FALSE,
15 #           times = 1)
16 # t_seti returns a matrix of indicies of rows that collectively
17 # comprise the training set and these rows are assigned to tr_set.
18 # This training set will be used to train the models.
19 tr_set<-wdbc.dt[t_seti,]
20 # The remaining 30% of rows comprise the test set and these can be
21 # obtained by a simple complement or wdbc.dt[-t_seti,].
22 ts_set<-wdbc.dt[-t_seti,]
23 # The next step is to convert all categorical variable diagnosis to
24 # numeric by assigning 'benign'=1 and 'malignant'=2.

```

```
25 tr_set$diagnosis<-c(benign=1,malignant=2)[tr_set$diagnosis]
26 ts_set$diagnosis<-c(benign=1,malignant=2)[ts_set$diagnosis]
27 # Convert resultant train data table to matrix using as.matrix()
28 tr_set<-as.matrix(tr_set)
29 # Now let us seperate dependant and independent variables.
30 # The second column is the independent variable and it is assigned
31 # to ytr_set.
32 ytr_set<-tr_set[,2]
33 # Columns 3 to 32 are the dependent variables and are is assigned
34 # to xtr_set.
35 xtr_set<-tr_set[,3:32]
36 # First we fit the Lasso Regression model (alpha=1) using
37 # cross-validation to the training data using the cv.glmnet()
38 # function.
39 fit.cv.lasso<-cv.glmnet(xtr_set,ytr_set,
40                         type.measure = 'auc',
41                         family='binomial')
42 # Next we fit the Ridge Regression model (alpha=0) in the
43 # same way.
44 fit.cv.ridge<-cv.glmnet(xtr_set,ytr_set,
45                         alpha=0,
46                         type.measure = 'auc',
47                         family='binomial')
48 # Next, we split the screen into 2 halves using par() with
49 # mfrow=c(1,2). Margin sizes can also be customized using
50 # 'mar' argument.
51 par(mfrow=c(1,2),
52      mar=c(4,4,5,2))
53 # Now we plot the Ridge and Lasso Regression objects.
54 plot(fit.cv.lasso,
55       main='Lasso Regression')
56 plot(fit.cv.ridge,
57       main='Ridge Regression')
```

Table 1: (1a)Lambda that maximises AUC

	Lambda
Lasso Regression	0.0158513
Ridge Regression	0.1042908



```

1 ## Answer in this chunk

1 # The lambda values that maximise AUC for both Regression models
2 # (lambda.min) are tabulated in a matrix.
3 res<-matrix(c(fit.cv.lasso$lambda.min,
4                 fit.cv.ridge$lambda.min),
5                 nrow=2,
6                 dimnames = list(c('Lasso Regression',
7                               'Ridge Regression'),
8                               c('Lambda'))))
9 # The matrix is returned using the kable() function.
10 kable(res,
11        caption = '(1a)Lambda that maximises AUC')

```

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** that contains the corresponding λ , AUC and model size.
- Use 3 significant figures for floating point values and comment on these results.
- *Note : The AUC values are stored in the field called `cvm`.*

```

1 # Store indices for lambda.min and lambda.1se for both Ridge
2 # and Lasso Regression respectively.
3 i1<-fit.cv.ridge$index[1]
4 i2<-fit.cv.ridge$index[2]
5 i3<-fit.cv.lasso$index[1]
6 i4<-fit.cv.lasso$index[2]
7 # Create a data table for respective lambda values, corresponding
8 # AUC and Model size
9 d<-data.table(
10   'Regression Type'=c('Ridge (lambda.min)',
11                       'Ridge (lambda.1se)',
12                       'Lasso (lambda.min)',
13                       'Lasso (lambda.1se)'),
14   'Lambda'=c(round(fit.cv.ridge$lambda[i1],3),
15             round(fit.cv.ridge$lambda[i2],3),
16             round(fit.cv.lasso$lambda[i3],3),
17             round(fit.cv.lasso$lambda[i4],3)),
18   'AUC'=c(round(fit.cv.ridge$cvm[i1],3),
19            round(fit.cv.ridge$cvm[i2],3),
20            round(fit.cv.lasso$cvm[i3],3),
21            round(fit.cv.lasso$cvm[i4],3)),
22   'Model Size'=c(fit.cv.ridge$nzero[i1],
23                  fit.cv.ridge$nzero[i2],
24                  fit.cv.lasso$nzero[i3],
25                  fit.cv.lasso$nzero[i4])
26 )
27 # Print data table using kable()
28 kable(d,caption = '(1b)Penalty Models Parameters')

1 # The first notable observation from the results is that the
2 # AUC in both Regression models are particularly bigger
3 # at lambda.min. This is actually because that is the value of
4 # lambda that maximizes AUC. The corresponding AUC that we see
5 # when lambda=lambda.min is the highest possible AUC you can get.

```

Table 2: (1b)Penalty Models Parameters

Regression Type	Lambda	AUC	Model Size
Ridge (lambda.min)	0.104	0.973	30
Ridge (lambda.1se)	341.505	0.964	30
Lasso (lambda.min)	0.016	0.976	7
Lasso (lambda.1se)	0.195	0.966	2

```

6 # Based on AUCs alone, Lasso at lambda.min appears to be the
7 # best model.
8 # Apart from this, the number of non-zero coefficients in Ridge
9 # Regression seem to be significantly bigger than Lasso Regression.
10 # This is because in Lasso Regression, some coefficients are shrunk
11 # to zero in the training process and thus gradually reducing the
12 # overall number of non-zero coefficients.

```

Problem 1.c (7 points)

- Perform both backward (we denote this as **model B**) and forward (**model S**) stepwise selection on the same training set derived in **problem 1.a**. Mute all the trace by setting **trace = FALSE**.
- Report the variables selected and their standardised regression coefficients in increasing order of the absolute value of their standardised regression coefficient.
- Discuss the results and how the different variables entering or leaving the model influenced the final result.
- **Note :** You can mute the warning by assigning **{r warning = FALSE}** for the chunk title

```

1 # First we create a data table with all required columns form the
2 # training set. The first column is unnecessary so we select all
3 # columns from the second one [2:32].
4 Tr_set<-data.table(tr_set[,2:32])
5 # Now we fit a full model(diagnosis~.) in which outcome is dependent
6 # on all 30 biomarkers using lm().
7 lin_mod<-lm(diagnosis~.,
8           data = Tr_set)
9 # Now we model stepwise backward selection. This can be done using
10 # the stepAIC() function with direction='backward'. We pass the
11 # full model (lin_mod) to it. The trace of the model explains every
12 # step in backward selection and for the moment, we will be skipping

```

Table 3: (1c)Model B: Variables Selected

x
radius
texture
perimeter
area
compactness
concavepoints
radius.stderr
concavity.stderr
radius.worst
area.worst
smoothness.worst
symmetry.worst
fractaldimension.worst

```

13 # that information by setting trace=FALSE.
14 sel.backward<-stepAIC(lin_mod,
15   direction = 'backward',
16   trace = FALSE
17 )
18 # Next, we report the variables selected. This can be done by
19 # passing the colnames() function on the new model. And since
20 # we are not interested in the outcome variable, we return
21 # all variables except that ([2:ncol(sel.backward$model)]).
22 # The result is returned using the kable() function.
23 kable(colnames(sel.backward$model)[2:ncol(sel.backward$model)],
24       caption = '(1c)Model B: Variables Selected')

1 # Finally, we report the coefficients of the model (stored in
2 # sel.backward$coefficients in increasing order of absolute values)
3 kable(
4   sel.backward$coefficients[order(abs(sel.backward$coefficients))],
5   caption = '(1c)Model B: Coefficients')

1 # Now we fit a null model (diagnosis~1) in which outcome is only
2 # dependent on the intercept of the model using lm().
3 fwd_model<-lm(diagnosis~1,
4   data=Tr_set)

```

Table 4: (1c)Model B: Coefficients

	x
area	-0.0004414
area.worst	-0.0008836
perimeter	-0.0059521
texture	0.0152970
radius	0.0481229
radius.worst	0.1540157
radius.stderr	0.3286845
symmetry.worst	0.7282311
concavity.stderr	1.2530096
(Intercept)	-1.5785245
smoothness.worst	2.7875483
compactness	-3.7735871
fractaldimension.worst	4.8277913
concavepoints	6.6934393

```

5 forward_select<-stepAIC(fwd_model,
6   scope = list(upper=lin_mod),
7   direction = 'forward',
8   trace=FALSE)
9 # Next, we report the variables selected. This can be done by passing
10 # the colnames() function on the new model. And since we are not
11 # interested in the outcome variable, we return all variables except
12 # that ([2:ncol(forward_select$model)]).
13 # The result is returned using the kable() function.
14 kable(
15 colnames(forward_select$model)[2:ncol(forward_select$model)],
16 caption = '(1c)Model S: Variables Selected')

1 # Finally, we report the coefficients of the model (stored in
2 # forward_select$coefficients in increasing order of absolute values)
3 kable(
4 forward_select$coefficients[order(abs(forward_select$coefficients))],
5 caption = '(1c)Model S: Coefficients')

1 # Backward selection starts with a full model and at each step
2 # eliminates a variable that has little or no significance to the
3 # prediction process. We see a small increase in the AIC value with

```

Table 5: (1c)Model S: Variables Selected

x
concavepoints.worst
radius.worst
area
texture
concavepoints
area.worst
perimeter.worst
radius.stderr
smoothness.worst
symmetry.worst
compactness
concavity.stderr
fractaldimension.worst
radius
perimeter

Table 6: (1c)Model S: Coefficients

	x
area	-0.0003773
area.worst	-0.0008454
perimeter.worst	-0.0048292
perimeter	-0.0049809
texture	0.0156029
radius	0.0448543
radius.worst	0.1713112
radius.stderr	0.3631839
symmetry.worst	0.7236488
concavepoints.worst	0.8301920
concavity.stderr	1.1258167
(Intercept)	-1.4836773
smoothness.worst	2.5027397
compactness	-3.4819636
fractaldimension.worst	4.4953955
concavepoints	5.6532651

Table 7: (1d)Pchisq-value

x
0.943856

```

4 # every step. Also expect a gradual decrease in prediction error
5 # with each step.
6 # Forward selection on the other hand starts with a null model.
7 # We start with trying to fit the best one predictor model that has
8 # lead to statistically significant reduction of AIC. The next step
9 # continues with all possible 2 predictor models and deciding the one
10 # that lead to significant decrease of AIC. This goes on and on
11 # until we get included required number of variables. Also expect a
12 # gradual decrease in prediction error with each step. trying to fit
13 # the best one predictor model that is statistically significant to
14 # outcome.

```

###Problem 1.d (3 points)

- Compare the goodness of fit of **model B** and **model S**
- Interpret and explain the results you obtained.
- Report the values using `kable()`.

```

1 # Here we carry out the p-chi square test to test the goodness
2 # of fit of models B and S.
3 # First, we calculate degrees of freedom which is simply the
4 # absolute difference in number of variables of both models.
5 dfreedom<-abs(ncol(sel.backward$model)-ncol(forward_select$model))
6 # We calculate p-chisquare at "dfreedom" degrees of freedom.
7 p<-pchisq(deviance(sel.backward)-deviance(forward_select),
8             df=dfreedom,
9             lower.tail = FALSE)
10 # Return p-value using kable().
11 kable(p,caption = '(1d)Pchisq-value')

1 kable(c(summary(sel.backward)$adj.r.square,
2 summary(forward_select)$adj.r.square),
3 caption = '(1d)Adjusted R Squared values Models: B and S')

```

Table 8: (1d) Adjusted R Squared values Models: B and S

x
0.7288784
0.7287561

```

1 ### Analysing Goodness of fit:
2 # Upon computing the pchisq-square value between models B and S,
3 # we observe that the value is very high (94.38%) than significance
4 # level alpha (0.05). This indicates that the models are not
5 # significantly different from each other and have more or less the
6 # same distribution.
7 # Furthermore, on comparing the adjusted r squared values, the
8 # backward selection model seems to have a slight edge over forward
9 # selection.

```

Problem 1.e (2 points)

- Plot the ROC curve of the trained model for both **model B** and **model S**. Display with clear title, label and legend.
- Report AUC values in 3 significant figures for both **model B** and **model S** using **kable()**.
- Discuss which model has a better performance.

```

1 # Plot an ROC curve to compare both models.
2 suppressMessages(invisible({
3   # First, create roc objects for both models (roc1
4   # and roc2 respectively).
5   roc1<-roc(Tr_set$diagnosis,
6             sel.backward$fitted.values)
7   roc2<-roc(Tr_set$diagnosis,
8             forward_select$fitted.values)
9   # Now we plot both roc curves with color of each plot specified
10  # using the 'col' argument.
11  plot(roc1,main='ROC Curves: Model B v Model S',col='red')
12  # The add=TRUE argument ensures that this plot is added to the
13  # first plot.
14  plot(roc2,add=TRUE,
15       col='purple')
16  # Finally, we create a legend for the combined plot using the

```

Table 9: (1e) Stepwise models and their AUCs

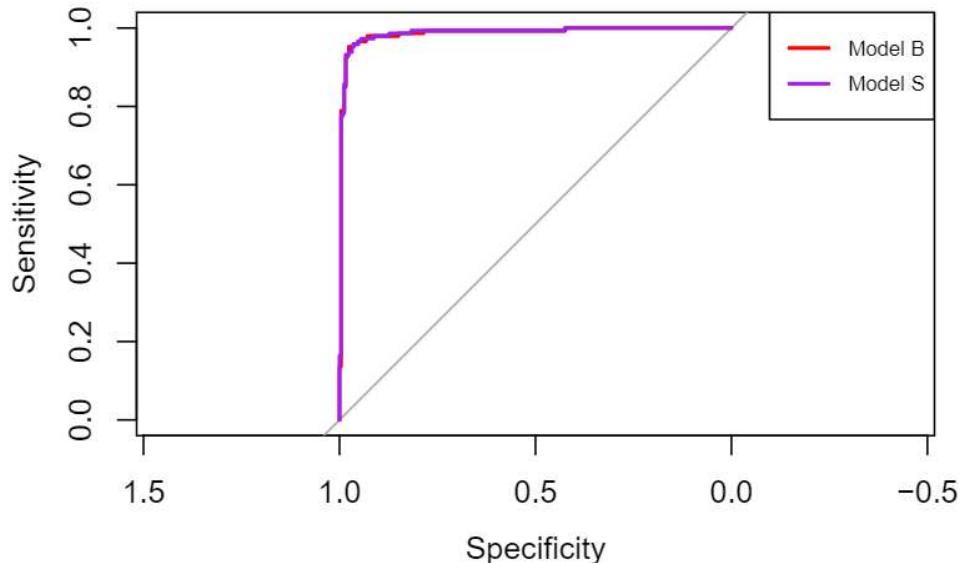
Model	AUC
Model B	0.987
Model S	0.987

```

17  # legend function.
18  legend('topright',
19    legend = c('Model B','Model S'),
20    col=c('red','purple'),
21    lty=c(1,1),lwd = c(2,2),cex=0.7)
22 }))

```

ROC Curves: Model B v Model S



```

1 # The final objective is to report the AUC of both curves. This
2 # value is stored in the 'auc' attribute of the roc object. The
3 # round() function is used to round the value to 3 decimal places.
4 results<-data.table('Model'=c('Model B','Model S'),
5   'AUC'=c(round(roc1$auc,3),round(roc2$auc,3)))
6 )
7 kable(results,
8   caption='(1e)Stepwise models and their AUCs')

```

```

1  ### Which model has better performance?
2  # The AUCs in both models appear to be almost the same with Model S
3  # being slightly better than Model B (0.9872 v 0.9869). Typically,
4  # we can expect both models in this case to perform almost the same.
5  # Since the number of variables under consideration is large (30
6  # variables), forward selection seems to be a much better alternative
7  # than backward selection. Thus, we see that forward selection
8  # slightly has the edge.

```

Problem 1.f (6 points)

- Use the four models to predict the outcome for the observations in the test set (use the λ 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.
- Display with clear title, label and legend.
- Compare the training AUCs obtained in problems 1.b and 1.e with the test AUCs and discuss the fit of the different models.

```

1  # Convert test data table to matrix using as.matrix()
2  ts_set<-as.matrix(ts_set)
3  # Next we separate the outcome variable and the predictor variables
4  # as yts_set and xts_set respectively.
5  # Columns 3 to 32 are the dependent variables and are assigned
6  # to xtr_set.
7  xts_set<-ts_set[,3:32]
8  # The second column is the independent variable and it is assigned
9  # to yts_set.
10 yts_set<-ts_set[,2]
11 # Now we make predictions on outcome using the predict() function
12 # from the glmnet class and thus we see the 'newx' argument.
13 # Lambda is set at 1 standard error for the penalized models
14 # as follows:
15 pred1<-predict(fit.cv.ridge,
16                  newx = xts_set,
17                  type='response',
18                  s='lambda.1se')
19 pred2<-predict(fit.cv.lasso,
20                  newx = xts_set,
21                  type='response',

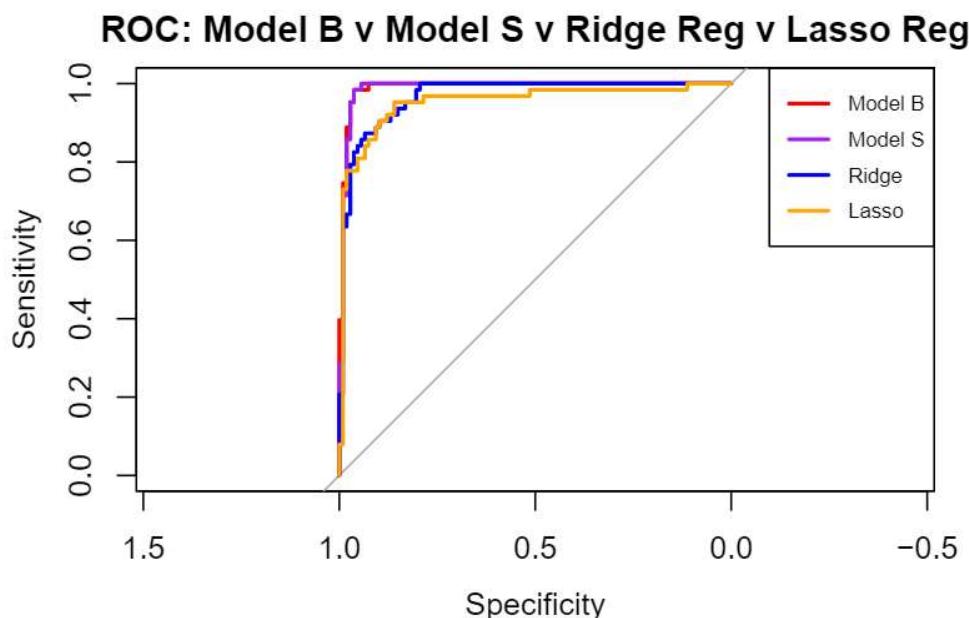
```

```

22           s='lambda.1se')
23   # Convert xts_set to data frame to carry out predictions for
24   # the stepwise selection models.
25   xts_set<-as.data.frame(xts_set)
26   # Next we make predictions for stepwise linear models. This time
27   # the predict() function has the 'newdata' argument.
28   # Predicted values for stepwise backward selection.
29   pred3<-predict(sel.backward,
30                     newdata = xts_set)
31   # Predicted values for stepwise forward selection.
32   pred4<-predict(forward_select,
33                     newdata = xts_set)
34   # 'suppressMessages' helps to suppress any unwanted messages that
35   # may be returned.
36   suppressMessages(invisible({
37     # Now we create roc() objects for all 4 models as follows:
38     roc1<-roc(yts_set,pred3)
39     roc2<-roc(yts_set,pred4)
40     roc3<-roc(yts_set,pred1)
41     roc4<-roc(yts_set,pred2)
42     # Now we plot all 4 ROC curve by passing the objects to the
43     # 'plot' function. To distinguish between plots, we give each plot
44     # a colour using the 'col' argument.
45     plot(roc1,
46           main='ROC: Model B v Model S v Ridge Reg v Lasso Reg',
47           col='red')
48     # add=TRUE adds the current plot to the first plot. We add the
49     # other plots in a similar way.
50     plot(roc2,add=TRUE,col='purple')
51     plot(roc3,add=TRUE,col='blue')
52     plot(roc4,add=TRUE,col='orange')
53     # Finally, we create a legend for the combined plot using the
54     # legend() function.
55     legend('topright',
56           legend = c('Model B','Model S','Ridge','Lasso'),
57           col=c('red','purple','blue','orange'),
58           lty=c(1,1),lwd = c(2,2),cex=0.7)
59
60 }))
```

Table 10: (1f)AUCs of models

Model	AUC
Model B	0.990
Model S	0.989
Ridge Reg	0.968
Lasso Reg	0.952



```

1 # The auc of each model can be found in the 'auc' attribute
2 # of the roc object. They are all tabulated into a data table
3 # as follows:
4 x<-data.table(
5   'Model'=c('Model B','Model S','Ridge Reg','Lasso Reg'),
6   'AUC' =c(round(roc1$auc,3),round(roc2$auc,3),
7             round(roc3$auc,3),round(roc4$auc,3)
8           ))
9 # Return the table using kable().
10 kable(x,caption = '(1f)AUCs of models')

1 # Based on the train test AUCs, it can be observed that Ridge
2 # Regression (0.968 > 0.964), Model B (0.990>0.987) and Model S
3 # (0.989>0.987) show higher test AUCs compared to training AUCs.

```

```
4  #
5  # But on the whole, the train and test AUCs seem to be almost the
6  # same and this shows that these models are ideal fits and generalize
7  # well.
8  #
9  # Finally, we have Lasso Regression which gives which gives a lower
10 # Test AUC compared to Train AUC ( $0.952 < 0.966$ ). Thus, this seems to
11 # generalize poorly compared to other models. Overall, Model B seems
12 # to have a better tradeoff between Train and Test AUC.
```

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 in file `GDM.raw.txt` into a data table named `gdm.dt`.
- Impute missing values in `gdm.dt` according to SNP-wise median allele count.
- Display first 10 rows and first 7 columns using `kable()`.

```

1 # In order to impute all null values to median we first create
2 # an impute.to.median() function that takes the data table
3 # and imputes missing values as median.
4 impute.to.median <- function(x) {
5     # only apply to numeric/integer columns
6     if (is.numeric(x) || is.integer(x)){
7         # find which values are missing
8         na.idx <- is.na(x)
9     # replace NAs with the median computed over the observed values
10        x[na.idx] <- median(x, na.rm = TRUE)
11    }
12    # return the vector with imputed values
13    return(x)
14 }
15 # Read the GDM data using fread()
16 gdm.dt<-fread('GDM.raw.txt')
17 # Store the column names with numeric values
18 numcols <- gdm.dt %>%
19   select_if(is.numeric) %>% colnames
20 # For all numeric columns, call impute.to.median()
21 gdm.dt <- gdm.dt %>%
22   .[, (numcols) := lapply(.SD,
23     #gdm.dt <- gdm.dt %>% copy %>%
24     .[, (numcols) := lapply(.SD,
25     # Return first 10 rows and first 7 columns using kable().
26   kable(head(gdm.dt,10)[,1:7],caption = '(2a) Imputed gdm.dt')
```

Table 11: (2a)Imputed gdm.dt

ID	sex	pheno	rs7513574_T	rs1627238_A	rs1171278_C	rs1137100_A
1	FALSE	0	1	0	0	2
2	FALSE	0	0	0	0	1
4	FALSE	1	2	1	1	1
5	FALSE	1	0	1	1	1
6	FALSE	1	0	1	1	1
7	FALSE	0	1	1	1	0
8	FALSE	0	0	0	0	1
12	FALSE	1	1	1	1	1
13	FALSE	1	2	0	0	2
18	FALSE	0	1	0	0	0

Problem 2.b (8 points)

- Write function `univ.glm.test()` where it takes 3 arguments, `x`, `y` and `order`.
- `x` is a data table of SNPs, `y` is a binary outcome vector, and `order` is a boolean which takes `false` as a default value.
- 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.

```

1 # We create a function univ.glm.test() that takes a data table
2 # of SNPs (x), outcomes (y) and order flag as arguments and applies
3 # Logistic Regression to all 176 SNPs.
4 univ.glm.test<-function(x,y,order=FALSE){
5   # Initialize a data frame 'res'.
6   res<-NULL
7   for(i in 1:ncol(x)){
8     # For all columns:
9     # Store i-th column
10    xi<-x[,i]
11    # Fit Logistic Regression model to i-th column or i-th SNP.
12    logistic_x<-glm(y~xi,family = binomial(link = 'logit'))
13    # For each SNP, we return SNP names, regression coefficients,
14    # odds ratios, standard errors and p-values.
15    # The whole set of values for each SNP is appended to res
16    # using rbind().
17    res<-rbind(res,

```

```

18   c(
19     colnames(x)[i],
20     logistic_x$coefficients[2],
21     exp(coef(logistic_x)[2]),
22     coef(summary(logistic_x))[-1, "Std. Error"],
23     coef(summary(logistic_x))[-1, "Pr(>|z|)"])
24   )
25 }
26 # Assign column names to resultant values
27 colnames(res)<-c('SNPName',
28   'Reg_Coefficient',
29   'Odds_Ratio',
30   'Std_Error',
31   'Pvalue')
32 # Convert data frame to data table
33 res<-data.table(res)
34 # Check if order=TRUE
35 if(order==TRUE){
36   # Order rows by increasing order of Pvalues.
37   res<-res[with(res,order(Pvalue))]
38 }
39 # Return res
40 return(res)
41 }
```

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`) and name the output data table as `gdm.as.dt`.
- Print the first 10 values of the output from `univ.glm.test()` using `kable()`.
- 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 using `kable()` from the GWAS.
- Report the 95% and 99% confidence intervals on the odds ratio using `kable()`.

```

1 # Store the SNP value columns into xtr in matrix form using
2 # as.matrix().
3 xtr<-as.matrix(gdm.dt[,4:ncol(gdm.dt)])
4 # Store outcome variable in ytr
```

```

5  ytr<-gdm.dt$pheno
6  # The outputs of the function univ.glm.test() are all stored in
7  # data table gdm.as.dt.
8  gdm.as.dt<-univ.glm.test(xtr,ytr)
9  # When we are looking for the SNP that is most strongly associated
10 # to increased risk of gestational diabetes, we are looking for the
11 # SNP that maximizes the chance of disease and thus the SNP with the highest odds ratio.
12 max_likely <-gdm.as.dt[Odds_Ratio==max(Odds_Ratio),
13                         SNPName]
14 # Returns "rs1423096_T"
15 # When we are looking for the SNP with most significant protective
16 # effect, we are looking for the SNP that minimizes the chance of
17 # disease and thus the SNP with the least odds ratio.
18 sig_pro <-gdm.as.dt[Odds_Ratio==min(Odds_Ratio),
19                      SNPName]
20 # Returns "rs11575839_C"
21 # We convert xtr from matrix to dataframe using data.frame(). This
22 # is done in order to pass it to the logistic regression model.
23 xtr<-data.frame(xtr)
24 # Fit a logistic regression model with outcome variable (pheno) and
25 # the strong risk associated SNP
26 logi_str_risk<-glm(ytr~rs1423096_T,
27                      data = xtr,
28                      family = binomial(link = 'logit'))
29 # Fit a logistic regression model with outcome variable (pheno) and
30 # the lowest risk associated SNP
31 logi_pro_effect<-glm(ytr~rs11575839_C,
32                      data = xtr,
33                      family = binomial(link = 'logit'))

1 # Return first 10 values of gdm.as.dt using head().
2 kable(head(gdm.as.dt,10),caption = 'gdm.as.dt')|>
3   kable_styling(full_width = F,
4                 position = "center",
5                 latex_options = "hold_position")

1 # Report summary statistics of the 2 focus SNPs
2 # (max_likely and sig_pro) using kable().
3 kable(coef(summary(logi_str_risk)),
4       caption = "(2c)SNP with Strong Risk Associated")

```

Table 12: gdm.as.dt

SNPName	Reg_Coefficient	Odds_Ratio	Std_Error	Pvalue
rs7513574_T	0.00215749439605489	1.00215982346177	0.105137211018395	0.983627959659055
rs1627238_A	0.114637864262422	1.1214672406608	0.113822428263324	0.313855918679796
rs1171278_C	0.121409444710005	1.12908711648435	0.11380728528091	0.286062822255494
rs1137100_A	0.060104751291512	1.06194778112087	0.1104238492224	0.586228503617821
rs2568958_A	0.149379896281938	1.16111400832769	0.123380016627863	0.225998873574341
rs1514175_A	0.0562296257329725	1.05784056326456	0.105235896141762	0.593120277525613
rs1555543_C	-0.0777805378818058	0.925167443472291	0.105783590355313	0.462169103637855
rs10923931_C	-0.214774420994388	0.806723399515719	0.18251171399461	0.239287076198348
rs516636_A	0.0552265983291012	1.05678005214045	0.123343524188903	0.654336351046674
rs574367_G	0.0590072860418808	1.06078296962115	0.124480094024075	0.635478588347025

Table 13: (2c)SNP with Strong Risk Associated

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	0.0824139	0.0734054	1.122722	0.2615556
rs1423096_T	0.6510641	0.3166547	2.056069	0.0397758

```

1 kable(coef(summary(logi_pro_effect)),
2   caption = "(2c)SNP with Most Significant Protective Effect")
3
4 # Initialize a dataframe to store 95% and 99% confidence intervals
5 # of the odds ratios of the 2 focus SNPs.
6 conf_ints<-NULL
7 # The first row will have 95% and 99% confidence limits of the
8 # strong risk SNP and this row is appended to conf_ints using rbind().
9 # The confidence limits are all rounded to 3 decimal places
10 # using round().
11 conf_ints<-rbind(conf_ints,
12   c('Strong Risk',
13     round(exp(confint(logi_str_risk)),3)[2,],
14     round(exp(confint(logi_str_risk,level=0.99)),3)[2,])

```

Table 14: (2c)SNP with Most Significant Protective Effect

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	0.1427220	0.0728282	1.959707	0.0500300
rs11575839_C	-0.6022542	0.3758156	-1.602526	0.1090394

```
12 ))
```

Waiting for profiling to be done...
 Waiting for profiling to be done...

```
1 # The second row will have 95% and 99% confidence limits of the
2 # protective effect SNP () and this row is appended to conf_ints
3 # using rbind(). The confidence limits are all rounded to 3
4 # decimal places using round().
5 conf_ints<-rbind(conf_ints,
6   c('Protective Effect',
7     round(exp(confint(logi_pro_effect)), 3)[2,],
8     round(exp(confint(logi_pro_effect,level=0.99)), 3)[2,]
9   ))
```

Waiting for profiling to be done...
 Waiting for profiling to be done...

```
1 # Assign column names to resultant dataframe.
2 colnames(conf_ints)<-c('SNP Type',
3   '95% CI:(2.5%)',
4   '95% CI:(97.5%)',
5   '99% CI:(0.5%)',
6   '99% CI:(99.5%)'
7 )
8 # Print conf_ints
9 kable(conf_ints,
10   caption = '(2c)Odds Ratio Confidence Intervals (by SNP')|>
11   kable_styling(full_width = F, position = "center",
12     latex_options = "hold_position")
```

Table 15: (2c)Odds Ratio Confidence Intervals (by SNP)

SNP Type	95% CI:(2.5%)	95% CI:(97.5%)	99% CI:(0.5%)	99% CI:(99.5%)
Strong Risk	1.05	3.67	0.874	4.567
Protective Effect	0.255	1.131	0.198	1.419

Problem 2.d (4 points)

- 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 $p\text{-value} < 10^{-4}$ (hit SNPs) report SNP name, effect allele, chromosome number, corresponding gene name and pos.
- Using `kable()`, report for each `snp.hit` the names of the genes that are within a 1Mb window from the SNP position on the chromosome.
- Note: That are genes that fall within +/- 1,000,000 positions using the pos column in the dataset.*

```

1 # Read file GDM.annot.txt into GDM.annot
2 GDM.annot<-fread('GDM.annot.txt')
3 # Now from a given SNP (Ex: s1234_X), we need to separately extract
4 # the id (s1234) and allele (X).
5 # snp id can be extracted from SNP by removing the last 2 characters
6 # ("s1234"_X) and that is extracted (substr(x,1,length(x)-2))
7 # as follows:
8 snp<-sapply(gdm.as.dt[,1],
9               function(i)
10              substr(i,1,nchar(i)-2))
11 # The last character in the snp string is the allele and that is
12 # extracted (substr(x,length(x),length(x))) as follows:
13 allele<-sapply(gdm.as.dt[,1],
14                  function(i)
15                  substr(i,nchar(i),nchar(i)))
16 # Create a copy of gdm.as.dt computed previously
17 gdm.as.dt_copy<-gdm.as.dt %>% copy
18 # We are going to replace SNPName column with two new columns
19 # and hence it is first set to null.
20 gdm.as.dt_copy[,SNPName:=NULL]
21 # The two new columns snp and allele are inserted into the
22 # data table
23 gdm.as.dt_copy[,`:=` (snp=snp,allele=allele)]
24 gdm.merged<-merge.data.table(gdm.as.dt_copy,
25                               GDM.annot,
26                               by.x='snp',
27                               by.y='snp',
28                               all = TRUE)

```

```

1 # The hit SNPs are first filtered out using as.numeric(Pvalue)<1e-4
2 # condition.
3 # Then we specify required columns: snp id, affected allele,
4 # chromosome number, gene and pos and return first 10 rows using
5 # kable as follows:
6 kable(head(gdm.merged[as.numeric(Pvalue)<1e-4,
7   c('snp','allele','chrom','gene','pos')],10),
8   caption = '2d:Hit SNPs') |>
9   kable_styling(full_width = F, position = "center",
10   latex_options = "hold_position")

```

Table 16: 2d:Hit SNPs

snp	allele	chrom	gene	pos
rs12243326	A	10	TCF7L2	114788815
rs2237897	T	11	KCNQ1	2858546

```

1 # Next for each hit SNP, we return genes (gene) that are within a
2 # 1Mb window from the SNP position (pos) on the chromosome. This
3 # process can be automated using sapply() on the hit SNP positions.
4 # Also note that there is a chance of repeated gene values that is
5 # eliminated by using unique().
6 # The result is stored in.snp.hit.genes.
7.snp.hit.genes<-sapply(gdm.merged[as.numeric(Pvalue)<1e-4, pos],
8   function(i)
9   unique(gdm.merged[as.numeric(Pvalue)<1e-4 & pos>=i-10^6 & pos<=i+10^6, gene]))
10 # Return the first 10 values using kable().
11 kable(head(snp.hit.genes,10),caption = '(2d)Genes (by hit SNP)')

```

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.

Table 17: (2d)Genes (by hit SNP)

x
TCF7L2
KCNQ1

- 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.
- Report odds ratio, 95% confidence interval and p-value using `kable()` for each score.

```

1 # First, let us store the rows that have Pvalue<10^-4 in SNPs_grs1.
2 SNPs_grs1<-gdm.as.dt[as.numeric(Pvalue)<1e-4]
3 # The SNPs that satisfy Pvalue<10^-4 in the 'SNPName' column
4 # of gdm.as.dt are all passed to SDcols. SDcols helps to select
5 # these columns alone and pass them to gdm_dt_snps1.
6 gdm_dt_snps1<- gdm.dt[, .SD,
7                         .SDcols = gdm.as.dt[as.numeric(Pvalue) < 1e-4]$SNPName]
8 # Check ordering of SNPs using the statement below:
9 # stopifnot(colnames(gdm.dt)[-1] == gdm.as.dt$SNPName)
10 # Now we do a matrix product of gdm_dt_snps1 and Regression
11 # coefficients stored in SNPs_grs1 to get weighted score:
12 weighted.score1 <- as.matrix(gdm_dt_snps1) %*%
13   as.numeric(SNPs_grs1$Reg_Coefficient)
14
15 # First, let us store the rows that have Pvalue<10^-3 in SNPs_grs2.
16 SNPs_grs2<-gdm.as.dt[as.numeric(Pvalue)<1e-3]
17 # The SNPs that satisfy Pvalue<10^-3 in the 'SNPName' column
18 # of gdm.as.dt are all passed to SDcols. SDcols helps to select
19 # these columns alone and pass them to gdm_dt_snps2.
20 gdm_dt_snps2 <- gdm.dt[, .SD,
21                         .SDcols = gdm.as.dt[as.numeric(Pvalue)<1e-3]$SNPName]
22 # Check ordering of SNPs using the statement below:
23 # stopifnot(colnames(gdm.dt)[-1] == gdm.as.dt$SNPName)
24 # Now we do a matrix product of gdm_dt_snps2 and Regression
25 # coefficients stored in SNPs_grs2 to get weighted score:
26 weighted.score2 <- as.matrix(gdm_dt_snps2) %*%
27   as.numeric(SNPs_grs2$Reg_Coefficient)
28
29 # First let us store the rows that have gene=='FTO' in SNPs_grs3.
30 SNPs_grs3<-gdm.merged[gene=='FTO']
31 # The SNPs that satisfy gene = 'FTO' in the 'SNPName' column
32 # of gdm.as.dt are all passed to SDcols. SDcols helps to select
33 # these columns alone and pass them to gdm_dt_snps3.
34 gdm_dt_snps3 <- gdm.dt[, .SD,
35                         .SDcols = gdm.merged[gene=='FTO',
36                                         paste(snp,allele,sep='_)]]
37 # Check ordering of SNPs using the statement below:

```

```

38 # stopifnot(colnames(gdm.dt)[-1] == gdm.merged$snp)
39 # Now we do a matrix product of gdm_dt_snps3 and Regression
40 # coefficients stored in SNPs_grs3 to get weighted score:
41 weighted.score3 <- as.matrix(gdm_dt_snps3) %*%
42   as.numeric(SNPs_grs3$Reg_Coefficient)
43 # Lets add all of these weighted scores to gdm.dt
44 gdm.dt<-gdm.dt[, `:=` (wgrsscore1=weighted.score1,
45                         wgrsscore2=weighted.score2,
46                         wgrsscore3=weighted.score3)]
47 # Create Logistic Regression models with 'pheno' being the outcome
48 # variable and weighted score 1 being predictor.
49 mod.weighted1<-glm(pheno~wgrsscore1,
50                     data = gdm.dt,
51                     family = binomial(link = 'logit'))
52 # Create Logistic Regression models with 'pheno' being the outcome
53 # variable and weighted score 2 being predictor.
54 mod.weighted2<-glm(pheno~wgrsscore2,
55                     data = gdm.dt,
56                     family = binomial(link = 'logit'))
57 # Create Logistic Regression models with 'pheno' being the outcome
58 # variable and weighted score 3 being predictor.
59 mod.weighted3<-glm(pheno~wgrsscore3,
60                     data = gdm.dt,
61                     family = binomial(link = 'logit'))
62 # Now we store the odds ratio, 95% Confidence interval and P-value
63 # of the 3 weighted models in data table 'result_table'.
64 result_table<-data.table(
65   'Weighted Score'=c(1,2,3),
66   'Odds Ratio'=c(
67     round(exp(coef(mod.weighted1)[2]),3),
68     round(exp(coef(mod.weighted2)[2]),3),
69     round(exp(coef(mod.weighted3)[2]),3)
70   ),
71   '95% CI:(2.5%)'=c(
72     round(confint(mod.weighted1)[2,1],3),
73     round(confint(mod.weighted2)[2,1],3),
74     round(confint(mod.weighted3)[2,1],3)
75   ),
76   '95% CI:(97.5%)'=c(
77     confint(mod.weighted1)[2,2],
78     confint(mod.weighted2)[2,2],

```

Table 18: (2e)Weighted Scores Model Results

Weighted Score	Odds Ratio	95% CI:(2.5%)	95% CI:(97.5%)	Pvalue
1	2.729	0.655	1.3638064	0.000
2	1.452	0.248	0.5014574	0.000
3	1.414	-0.200	0.8971549	0.215

```

79     confint(mod.weighted3)[2,2]
80             ),
81 'Pvalue'=c(
82     round(coef(summary(mod.weighted1))[-1, "Pr(>|z|)"],3),
83     round(coef(summary(mod.weighted2))[-1, "Pr(>|z|)"],3),
84     round(coef(summary(mod.weighted3))[-1, "Pr(>|z|)"],3)
85   )
86

```

Waiting for profiling to be done...
 Waiting for profiling to be done...

```

1 # Print result_table using kable()
2 kable(result_table,
3       caption = '(2e)Weighted Scores Model Results')

```

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*).

```

1 # Read the test set using fread()
2 gdm.test<-fread('GDM.test.txt')
3 # Find the SNPs in gdm.merged[as.numeric(Pvalue)<1e-4] that are
4 # in gdm.test and store in 'x1'.
5 x1<-intersect(colnames(gdm.test)[4:179],
6                 gdm.merged[as.numeric(Pvalue)<1e-4,snp])
7 # Filter out the rows using as.numeric(Pvalue)<1e-4 and 'x1'
8 SNPs_grs4<-gdm.merged[as.numeric(Pvalue)<1e-4 & SNP %in% x1]
9 # Extract the columns from 'x1'
10 gdm_dt_snps4<- gdm.test[, .SD,.SDcols = x1]
11 # Now we do a matrix product of gdm_dt_snps4 and Regression
12 # coefficients stored in SNPs_grs4 to get weighted score:
13 weighted.score4 <- as.matrix(gdm_dt_snps4) %*%
14   as.numeric(SNPs_grs4$Reg_Coefficient)
15
16 # Find the SNPs in gdm.merged[as.numeric(Pvalue)<1e-3] that are in
17 # gdm.test and store in 'x2'.
18 x2<-intersect(colnames(gdm.test)[4:179],
19                 gdm.merged[as.numeric(Pvalue)<1e-3,snp])
20 # Filter out the rows using as.numeric(Pvalue)<1e-3 and 'x2'
21 SNPs_grs5<-gdm.merged[as.numeric(Pvalue)<1e-3 & SNP %in% x2]
22 # Extract the columns from 'x2'
23 gdm_dt_snps5<- gdm.test[, .SD,.SDcols = x2]
24 # Now we do a matrix product of gdm_dt_snps5 and Regression
25 # coefficients stored in SNPs_grs5 to get weighted score:
26 weighted.score5 <- as.matrix(gdm_dt_snps5) %*%
27   as.numeric(SNPs_grs5$Reg_Coefficient)
28
29 # Find the SNPs in gdm.merged[gene=='FTO'] that are in gdm.test
30 # and store in 'x3'.
31 x3<-intersect(colnames(gdm.test)[4:179],
32                 gdm.merged[gene=='FTO',snp])
33 # Filter out the rows using gene=='FTO' and 'x3'
34 SNPs_grs6<-gdm.merged[ gene=='FTO' & SNP %in% x3]
35 # Extract the columns from 'x3'
36 gdm_dt_snps6<- gdm.test[, .SD,.SDcols = x3]
37 # Now we do a matrix product of gdm_dt_snps6 and Regression
38 # coefficients stored in SNPs_grs6 to get weighted score:
39 weighted.score6 <- as.matrix(gdm_dt_snps6) %*%
40   as.numeric(SNPs_grs6$Reg_Coefficient)
41 # Lets add all of these weighted scores to gdm.test

```

Table 19: (2g)Weighted Models

Weighted Model	Log-likelihood
1	-407.774
2	-1074.865
3	-1036.955

```

42 gdm.test<-gdm.test[, `:=` (wgrsscore1=weighted.score4,
43                               wgrsscore2=weighted.score5,
44                               wgrsscore3=weighted.score6)]

```

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 and present them using **kable()**

```

1 # Use the previous models to make predictions for the 3 genetic
2 # risk scores using the test data.
3 pred1<-predict(mod.weighted1,
4                  newdata=gdm.test)
5 pred2<-predict(mod.weighted2,
6                  newdata=gdm.test)
7 pred3<-predict(mod.weighted3,
8                  newdata=gdm.test)
9 # Next, calculate combined Log Likelihood for all 3 predictors. The
10 # dbinom() function is used for this purpose. We are interested in a
11 # single trial and hence size=1 and log=TRUE ensures that the
12 # probabilities
13 # are passed as# as
14 L<-c(sum(dbinom(gdm.dt$pheno,prob=pred1,size=1,log = TRUE),
15        na.rm = TRUE),
16 sum(dbinom(gdm.dt$pheno,prob=pred2,size=1,log = TRUE),
17      na.rm = TRUE),
18 sum(dbinom(gdm.dt$pheno,prob=pred3,size=1,log = TRUE)))
19 # Print data table using kable().
20 kable(data.table('Weighted Model'=c(1,2,3),'Log-likelihood'=L),
21       caption = '(2g)Weighted Models')

```

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*)
- 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 using `kable()`.

```

1 # Read summary statistics from GDM.study2.txt using fread()
2 gdm.study<-fread('GDM.study2.txt')
3 gdm.effect<- gdm.merged %>% copy
4 gdm.effect<-gdm.effect[with(gdm.effect,order(Pvalue))]
5 # Harmonise the variables of study
6 gdm.study<-gdm.study[snp %in% gdm.effect$snp]
7 gdm.effect<-gdm.effect[snp %in% gdm.study$snp]

1 gwas1 <- gdm.study[order(snp)]
2 gwas2 <- gdm.effect[order(snp)]
3 # all.equal(gdm.effect$snp,gdm.study$snp): TRUE
4
5 both.ok <-gwas2[as.numeric(Pvalue)<1e-4,allele] == gwas1$effect.allele
6 flipped <-gwas2[as.numeric(Pvalue)<1e-4,allele] == gwas1$other.allele
7
8 kable(
9   table(both.ok, flipped)
10  ,caption = "(2h)SNPs with Alleles Matching") |>
11  kable_styling(full_width = F, position = "center",
12    latex_options = "hold_position")

```

Table 20: (2h)SNPs with Alleles Matching

	FALSE	TRUE
FALSE	81	44
TRUE	51	0

```

1 # Initialize effect sizes
2 Beta1 <- as.numeric(gdm.effect$Reg_Coefficient)
3 Beta2 <- as.numeric(gdm.study$beta)
4 Beta2[flipped] <- -Beta2[flipped]
5 # Calculate inverse variance weights

```

```

6 wgwas1 <- 1 / as.numeric(gdm.effect$Std_Error)^2
7 wgwas2 <- 1 / gdm.study$se^2
8 # Print weights of first GWAS
9 kable(t(head(wgwas1)),
10       caption = "(2h)Weight of 1st GWAS (wgwas1)") |>
11     kable_styling(full_width = F, position = "center",
12                   latex_options = "hold_position")

```

Table 21: (2h)Weight of 1st GWAS (wgwas1)

81.38295	55.17323	52.35296	55.4531	49.45205	29.02107
----------	----------	----------	---------	----------	----------

```

1 # Print weights of second GWAS
2 kable(t(head(wgwas2)),
3       caption = "(2h)Weight of 2nd GWAS (wgwas2)") |>
4     kable_styling(full_width = F, position = "center", latex_options = "hold_position")

```

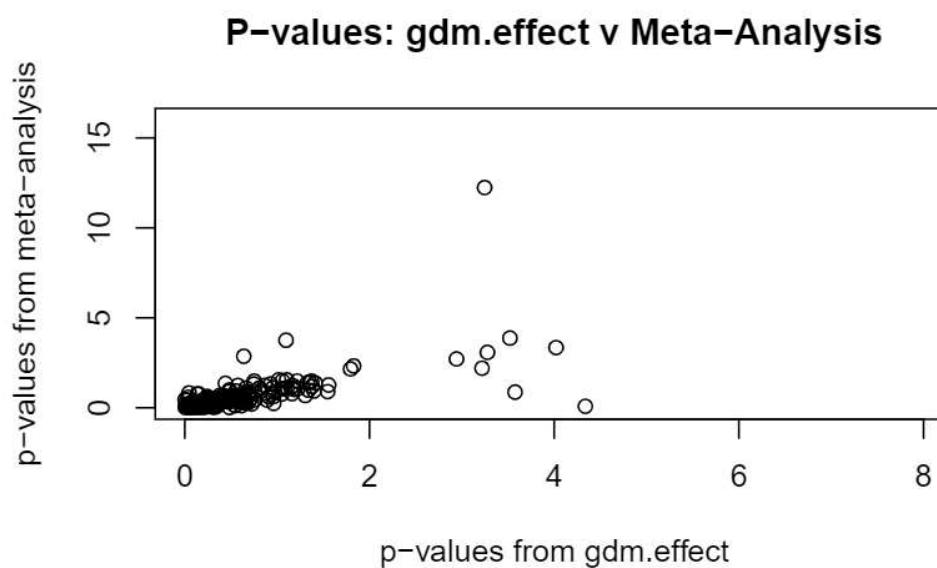
Table 22: (2h)Weight of 2nd GWAS (wgwas2)

34.82747	8.361486	1.397415	38.53447	11.27683	9.14147
----------	----------	----------	----------	----------	---------

```

1 # Upon observation, Study 1 looks more empowered
2
3 # Meta-analysis effect size (beta.ma) is a weighted sum of the
4 # effect sizes.
5 beta.ma <- (wgwas1 * Beta1 + wgwas2 * Beta2) /
6   (wgwas1 + wgwas2)
7 # Meta Analysis Standard Error
8 se.ma <- sqrt(1 / (wgwas1 + wgwas2))
9 # Meta Analysis P-values
10 pval.ma <- 2 * pnorm(abs(beta.ma / se.ma), lower.tail = FALSE)
11 # Plot P-values from 1st study v P-values from meta analysis
12 plot(-log10(as.numeric(gdm.effect$Pvalue)), -log10(pval.ma),
13       xlim = c(0, 8), ylim = c(0, 16),
14       xlab = "p-values from gdm.effect",
15       ylab = "p-values from meta-analysis",
16       main = "P-values: gdm.effect v Meta-Analysis"
17     )

```



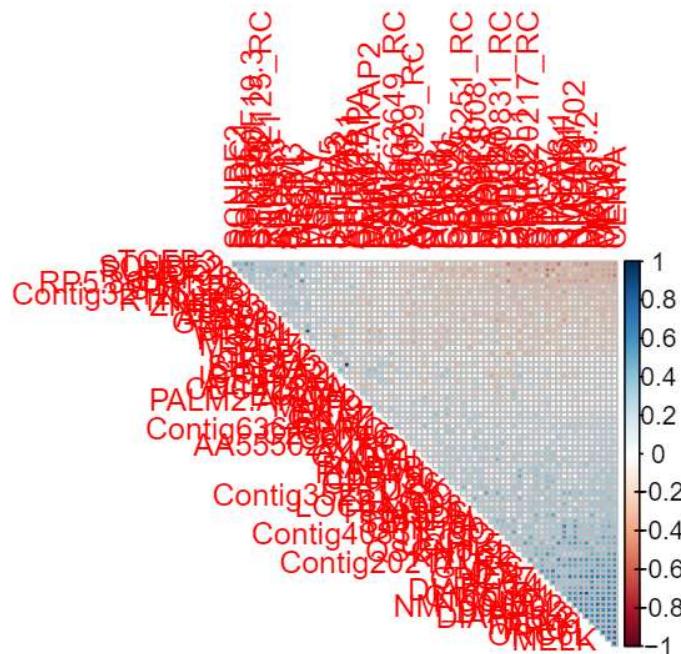
Problem 3 (33 points)

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 correlation matrix between the gene expression variables, and display it so that a block structure is highlighted using the `corrplot` package.
- Discuss what you observe.
- Identify the unique pairs of (distinct) variables that have correlation coefficient greater than 0.80 in absolute value and report their correlation coefficients.

```
1 # Read the data file using fread()
2 nki.dt<-fread('nki.csv')
3 # Create a new data table genes that contains all gene expression
4 # columns of nki.dt
5 genes<-data.table(nki.dt[,7:76])
6 # Store the list of names of numeric columns
7 numcols <- sapply(genes, is.numeric)
8 # Subset of numeric columns
9 cor.genes <- genes[, ..numcols] %>% #subset of numeric columns
10   cor(use="pairwise.complete")
11 # Plot a Correlation plot of the 70 gene expressions
12 corrplot(cor.genes, method = 'square',
13           order = 'FPC', type = 'upper', diag = FALSE)
```



```

1  #### What we observe?
2  # The method that we use to compute the coreelation coefficient
3  # is the default one (Pearson). Since  $\text{Corr}(A,B) = \text{Corr}(B,A)$ ,
4  # we can skip one triangle of values in the matrix. (It should
5  # also be noted that the correlation of a variable with itself
6  # is 1 and thus the diagonals are all blues based on the color
7  # marker to the right.)
8  # In our case, we will be skipping the lower triangle using
9  # type= 'upper'. The positive correlations are displayed in blue
10 # while the negative correlations are displayed red. As we specified
11 # method = 'square', each cell in the matrix will have a sqaure in it,
12 # the size of which is dependent on the corresponding correlation
13 # coefficient.

1  #### Identify variable pairs with cor > 0.8
2  K <- which(abs(cor.genes)>0.8 & row(cor.genes)<col(cor.genes),
3           arr.ind=TRUE)
4  # Reconstruct names from positions
5  Corr80 <- matrix(colnames(cor.genes)[K], ncol=2)
6  # Bind the columns into Corr80
7  Corr80<-cbind(Corr80,cor.genes[K])

```

Table 23: (3a) Variable pairs with Correlation Coefficient > 0.8

Var1	Var2	Corr_Value
DIAPH3	DIAPH3.1	0.803136819514865
DIAPH3	DIAPH3.2	0.833859118546205
DIAPH3.1	DIAPH3.2	0.88687406089682
PECI	PECI.1	0.869783649528927
IGFBP5	IGFBP5.1	0.97750296924928
NUSAP1	PRC1	0.829835551228507
PRC1	CENPA	0.817542383966119

```

8 # Assign column names
9 colnames(Corr80)<-c('Var1','Var2','Corr_Value')
10 # Print the pairs with corr(A,B)>0.8
11 kable(Corr80,
12   caption = '(3a)Variable pairs with Correlation Coefficient > 0.8')
```

Problem 3.b (8 points)

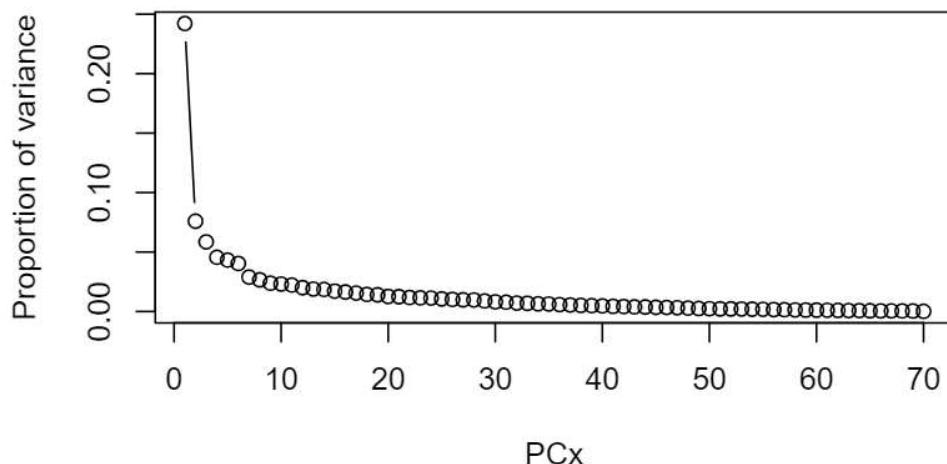
- Perform PCA analysis (only over the columns containing gene expressions) in order to derive a patient-wise summary of all gene expressions (dimensionality reduction).
- Decide which components to keep and justify your decision.
- Test if those principal components are associated with the outcome in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

```

1 # Apply the prcomp() function which does required PCA on 'genes'
2 pca.v<-prcomp(genes,center = T,scale. = T)
3 # Fraction of variance covered by each variable
4 perc.expl <- pca.v$sdev^2 / sum(pca.v$sdev^2)
5
6 # Plot a scree plot to look and decide on which principle components
7 # to keep
8 plot(perc.expl, main="Scree plot",
9       type = 'b',
10       xlab='PCx',
11       ylab='Proportion of variance'
12     )
```

Table 24: (3b)Proportion of variance covered

x
0.9039862

Scree plot

```

1 # Proportion of variance covered
2 kable(sum(perc.expl[1:33]),
3       caption = '(3b)Proportion of variance covered')

1 # Around 90.4%

1 #### Which principle components to keep?
2 # Upon looking at the scree plot, the curve flattens after 40.
3 # This means that the first 40 components cover a major portion
4 # of the variance (around 94.3%). So it would suffice to choose
5 # the first 30 to 40 components. We set our threshold to components
6 # that cover atleast 90% of the variance. The first 33 components
7 # are just enough to satisfy this criteria (around 90.4%). And thus,
8 # that is the new dimensional space. We have cut down 70 components
9 # to just 33.

```

```

1 # Create a copy of nki.dt
2 nki.dt.new<-nki.dt %>% copy
3 # Convert the categorical variables to numeric.
4 nki.dt.new$EstrogenReceptor<-c(Positive=1,
5                               Negative=2)[nki.dt.new$EstrogenReceptor]
6 nki.dt.new$Grade<-c(Intermediate=1,
7                      `Well diff`=2,
8                      `Poorly diff`=3)[nki.dt.new$Grade]
9 # Adjusted Model: Adjusted with Age, EstrogenReceptor and Grade as covariates.
10 nki.dt.new.lr<-glm(Event~Age+EstrogenReceptor+Grade,
11                     data = nki.dt.new,
12                     family = binomial(link = 'logit'))
13 # View the required coefficients in the summary of the model.
14 kable(coef(summary(nki.dt.new.lr)),
15       caption = 'Adjusted Model')|>
16   kable_styling(full_width = F, position = "center",
17                 latex_options = "hold_position")

```

Table 25: Adjusted Model

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	1.2884817	1.6465208	0.7825480	0.4338926
Age	-0.0676552	0.0344309	-1.9649576	0.0494191
EstrogenReceptor	0.6209388	0.4728793	1.3131021	0.1891485
Grade	0.1268021	0.2282015	0.5556586	0.5784443

```

1 # The unadjusted models are simply the models obtained by
2 # modelling the relationship between Event and each Principle
3 # Component separately (For each PCx, Event ~ PCx).
4 # Conveniently, the univ.glm.test() from 2b, does this for us.
5 Unadj.mod<-univ.glm.test(pca.v$x[,1:33],
6                           nki.dt.new$Event,
7                           order=TRUE)
8 # Let's change the column name 'SNPName' to 'Principle_Componentx'
9 # for better understanding
10 colnames(Unadj.mod)[1] <- "Principle_Componentx"
11 # Now let us sort out the models in which the Principle
12 # Components are statistically significant in the prediction
13 # process. This models have their Pvalue less than level of
14 # significance (0.05).
15

```

```

16 # Unadj.mod[as.numeric(Pvalue)<0.05, Principle_Componentx]:
17 # PC19,PC3,PC1,PC23,PC11,PC14

1 ### Adjusted and Non-Adjusted models:
2 # If we have a look at the adjusted model, only the age column seems
3 # to be statistical significant in the whole prediction process. This
4 # can be concluded as the other 2 variables, EstrogenReceptor and
5 # Grade have much higher P values.
6
7 # In case of the Non-adjusted models, there are only 6 models in
8 # which the predictor variables are of statistical significance.
9 # The major difference that we need to observe is the in case of the
10 # adjusted model, the results are based on an adjusted combination
11 # of more than one variable (atleast one covariate in addition)
12 # (Age, EstrogenReceptor and Grade variables in our case). Whereas,
13 # in case of Non-adjusted models, there are no additional covariates
14 # and the modelling is only based on 1 variable.

```

Problem 3.c (8 points)

- Use PCA plots to compare the main drivers with the correlation structure observed in problem 3.a.
- Examine how well the dataset may explain your outcome.
- Discuss your findings in full details and suggest any further steps if needed.

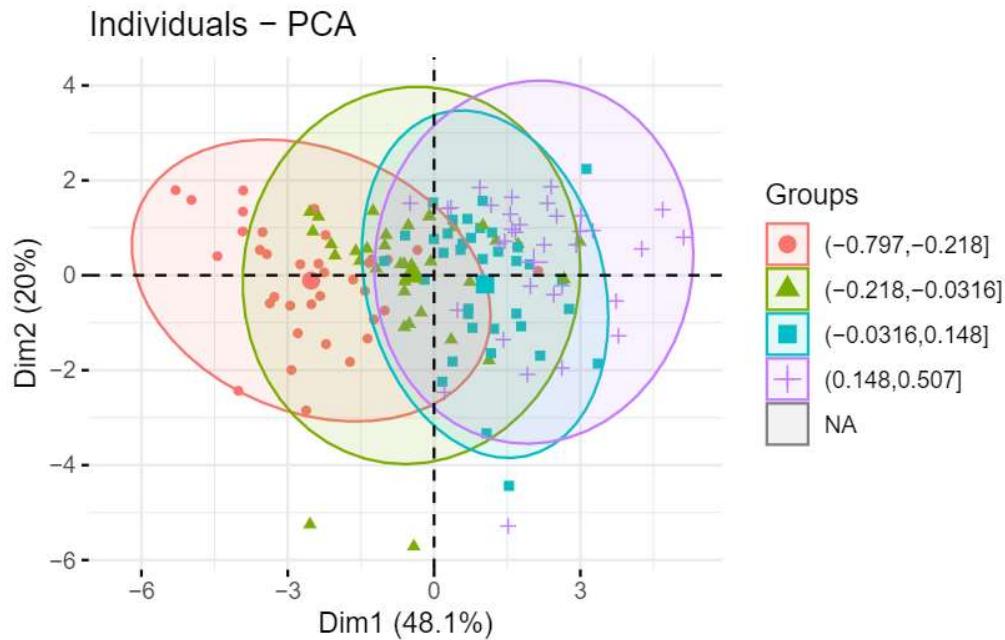
```

1 # Retrieve column names of the main drivers.
2 x<-c(unique(colnames(cor.genes)[K]))
3 # Apply pca on those drivers
4 pca.drivers<-prcomp(nki.dt[,..x],
5                      center = T,
6                      scale. = T)
7 # Now let us return all pca plots
8 # Graph of individuals
9 fviz_pca_ind(pca.drivers, geom = 'point',
10               habillage = cut(nki.dt$NUSAP1,
11                               #colour by NUSAP1
12                               quantile(nki.dt$NUSAP1)),
13               addEllipses = T)

```

Too few points to calculate an ellipse

Warning: Removed 1 rows containing missing values (`geom_point()`).
Removed 1 rows containing missing values (`geom_point()`).



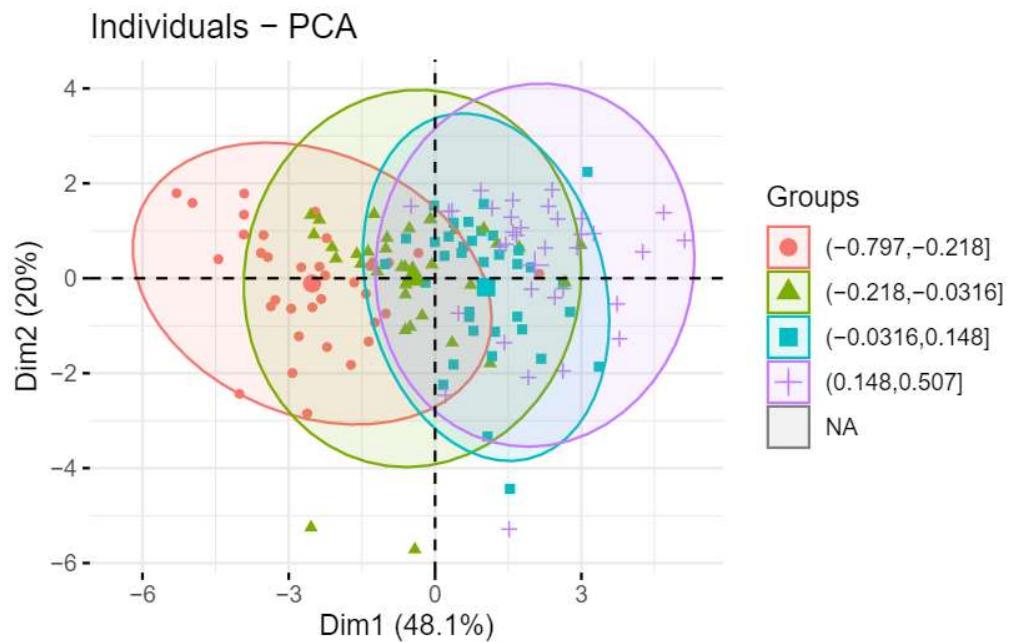
```

1 # Graph of individuals with data points projected along
2 # the second and third Principle Components.
3 fviz_pca_ind(pca.drivers, geom = 'point',
4               habillage = cut(nki.dt$NUSAP1, axes=c(2,3),
5               #colour by NUSAP1
6               quantile(nki.dt$NUSAP1)),
7               addEllipses = T)

```

Too few points to calculate an ellipse

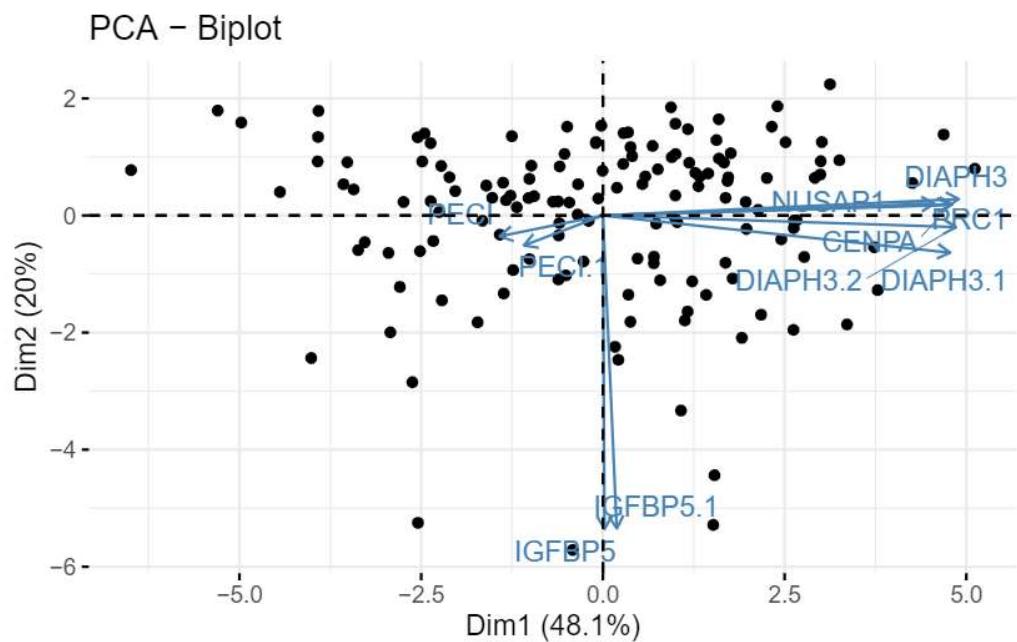
Warning: Removed 1 rows containing missing values (`geom_point()`).
Removed 1 rows containing missing values (`geom_point()`).



```

1 # Biplot
2 fviz_pca_biplot(pca.drivers, geom='point', repel = T)

```



```

1  ### Observations:
2  # From the first 3 plots we observe that the ellipses overlap with
3  # each other. The centroids (shown with much bigger data markers)
4  # are also separated in the right order along the x-axis although
5  # this separation is not seen along the y-axis.
6
7  # In the 2nd plot, we can look at the projections of variables drawn
8  # on a 2d plane. Upon looking at the extreme values (Values far
9  # left or far right of the axes), along Y-axis, we can observe that
10 # IGFBP5 and IGFBP5.1 have been assigned the highest magnitude by
11 # Principle Component 2. Similarly, along X-axis, we can observe that
12 # PECI and DIAPH3, DIAPH3.1 and PRC1 have been assigned the highest
13 # weights by Principle Component 1. This can also be observed in
14 # the correlation plot we did in 3a.
15
16 # From 3b, we know that only 6 variables (Out of 70) are associated
17 # with outcome (Event). This sums up the need to eliminate variables
18 # that are unassociated with outcome. An other idea might be to use
19 # the box plot for outlier detection.

```

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 with several experiments and evidences.

```

1  PCx<-data.frame(pca.v$x[,1:33])
2  # Create a Null Model with just the intercept as independent variable
3  Null.mod<-lm(nki.dt$Event~1,data=PCx)
4  # Create a Full Model with all independent variables
5  Full.mod<-lm(nki.dt$Event~,data=PCx)
6  # Now we use these models to fit Stepwise selection models
7  # Stepwise backward selection model
8  backward.sel<-stepAIC(Full.mod,
9                      direction = 'backward',
10                     scope = formula(Full.mod),
11                     trace = FALSE)
12 # Predictions based on above model
13 pred_1<-predict(backward.sel,newdata= PCx)
14

```

```

15 # Stepwise forward selection model
16 forward.sel<- stepAIC(Null.mod,
17                         direction = 'forward',
18                         scope = formula(Full.mod),
19                         trace = FALSE)
20 # Predictions based on above model
21 pred_2<-predict(forward.sel,newx= PCx)
22
23 Lasso.reg<-cv.glmnet(pca.v$x[,1:33],nki.dt$Event,
24                       type.measure = 'auc',
25                       family='binomial')
26 # Predictions based on above model
27 pred_3<-predict(Lasso.reg,newx=pca.v$x[,1:33])
28
29 # Ridge Regression
30 Ridge.reg<-cv.glmnet(pca.v$x[,1:33],nki.dt$Event,
31                       alpha=0,
32                       type.measure = 'auc',
33                       family='binomial')
34 pred_4<-predict(Ridge.reg,newx= pca.v$x[,1:33])
35
36 # Logistic Regression: Most accurate model
37 log_reg<-glm(nki.dt$Event~.,
38               data = PCx,
39               family = binomial(link='logit'))
40 pred_5<-predict(log_reg,newdata= PCx)

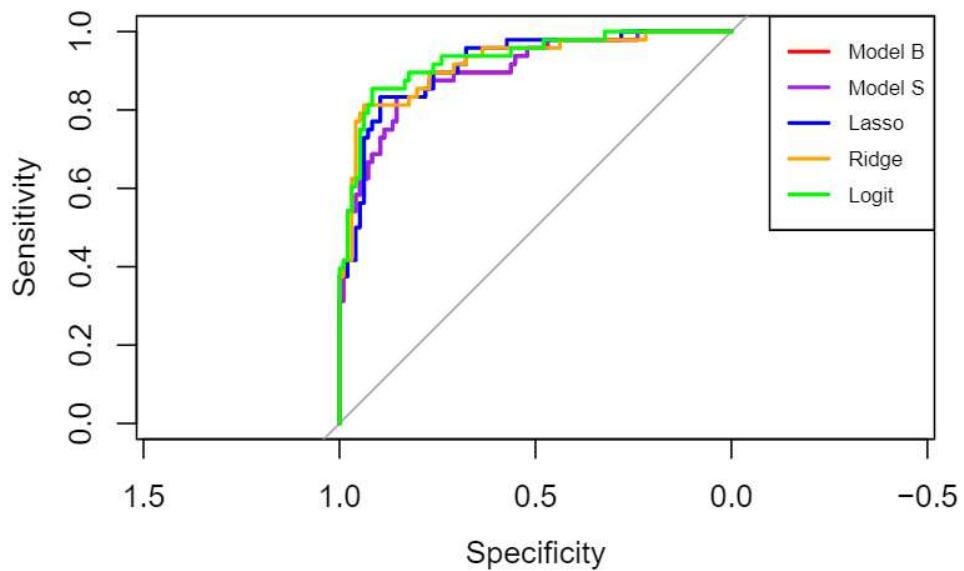
1 suppressMessages(invisible({
2   # Now we create roc() objects for all 4 models as follows:
3   roc1<-roc(nki.dt$Event,pred_1)
4   roc2<-roc(nki.dt$Event,pred_2)
5   roc3<-roc(nki.dt$Event,pred_3)
6   roc4<-roc(nki.dt$Event,pred_4)
7   roc5<-roc(nki.dt$Event,pred_5)
8   # Now we plot all 5 ROC curve by passing the objects to the
9   # 'plot' function. To distinguish between plots, we give each plot
10  # a colour using the 'col' argument.
11  plot(roc1,main='ROC: Model Performances',col='red')
12  # add=TRUE adds the current plot to the first plot. We add the other
13  # plots in a similar way.
14  plot(roc2,add=TRUE,col='purple')

```

```

15 plot(roc3,add=TRUE,col='blue')
16 plot(roc4,add=TRUE,col='orange')
17 plot(roc5,add=TRUE,col='green')
18 # Finally, we create a legend for the combined plot using the
19 # legend() function.
20 legend('topright',
21       legend = c('Model B','Model S','Lasso','Ridge','Logit'),
22       col=c('red','purple','blue','orange','green'),
23       lty=c(1,1),lwd = c(2,2),cex=0.7)
24
25 }))

```

ROC: Model Performances

```

1 # From observing the set of AOC curves from the above plot,
2 # Logistic Regression is the model that maximises AUC
3 # (AUC=0.9286)

```