

Assignment 1

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

March 9, 2023

Due on Thursday, 9th of March 2023, 5:00pm

! Pay Attention

The assignment is marked out of 100 points, and will contribute to **20%** of your final mark. The aim of this assignment is to produce a precise report in biomedical studies with the help of statistics and machine learning. Please complete this assignment using **Quarto/Rmarkdown** file and render/knit this document only in PDF format and submit using the **gradescope link on Learn**. 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.

Clear and reusable code will be rewarded. Codes without proper indentation, choice of variable identifiers, **comments**, error checking, 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 with `kable()` and `kable_styling()` or using [Markdown syntax](#). All plots must be displayed with clear title, label and legend.

Problem 1 (25 points)

Files `longegfr1.csv` and `longegfr2.csv` (available on Assessment > Assignment 1) contain information regarding a longitudinal dataset containing records on 250 patients. For each

subject, eGFR (estimated glomerular filtration rate, a measure of kidney function) was collected at irregularly spaced time points: variable `fu.years` contains the follow-up time (that is, the distance from baseline to the date when each eGFR measurement was taken, expressed in years).

Problem 1.a (4 points)

- Convert the files to data table format and merge in an appropriate way into a single data table.
- Order the observations according to subject identifier and follow-up time.
- Print first 10 values of the new dataset using `head()`.

```

1 # s2319494
2 # First let us read the csv files using the fread() function.
3 lgfr1<-fread('longegfr1.csv')
4 lgfr2<-fread('longegfr2.csv')
5 # Now we convert the data in lgfr1, lgfr2 into data tables using
6 # the setDT() function.
7 lgfr1<-setDT(lgfr1)
8 lgfr2<-setDT(lgfr2)
9 # The next step is to merge the data tables. In order to do so, we need
10 # common columns or keys. On observing, id and fu.years are the keys. The
11 # merge.data.table() function is then used to merge the data tables lgfr1
12 # and lgfr2.
13 lgfr<-merge.data.table(lgfr1,lgfr2,by.x = c('id','fu.years'),
14                         by.y = c('ID','fu.years'),all = TRUE)
15
16 # The next step is to order the entire table by id and follow-up time. This
17 # can be done conveniently using with() and thus specifying the columns to
18 # order by using order(). The results are stored in lgfr12
19 lgfr12<-lgfr[with(lgfr,order(id,fu.years))]
20 # View first 10 rows using head().
21 head(lgfr12,10)

```

	<code>id</code>	<code>fu.years</code>	<code>sex</code>	<code>baseline.age</code>	<code>egfr</code>
1:	1	0.0000	0	65.5	76.48
2:	1	0.1533	0	65.5	47.36
3:	1	0.6899	0	65.5	94.87
4:	1	1.1882	0	65.5	52.12
5:	1	1.8398	0	65.5	91.91
6:	1	2.2806	0	65.5	76.52

```

7: 1 3.3895 0      65.5 46.79
8: 1 3.7563 0      65.5 35.56
9: 1 4.5229 0      65.5 28.41
10: 1 5.3607 0     65.5 20.85

```

Problem 1.b (6 points)

- Compute the average eGFR and length of follow-up for each patient.
- Print first 10 values of the new dataset using `head()`.
- Tabulate the number of patients with average eGFR in the following ranges: $(0, 15]$, $(15, 30]$, $(30, 60]$, $(60, 90]$, $(90, \max(\text{eGFR}))$.
- Count and report the number of patients with missing average eGFR.

```

1 # Our objective now is to compute average or mean of eGFR and total length
2 # of follow up time (fu.years) by patient.
3 # We do so by computing two new columns avg_eGFR and sum(fu.years) in the process.
4 # And on specifying by='id', R does the computation by patient. The result is
5 # stored in Avg_eGFR_fu.
6 Avg_eGFR_fu<-lgfr12[,(avg_eGFR=mean(egfr),length_of_followup=sum(fu.years)),by='id']
7
8 # Next we print first 10 values using head()
9 head(Avg_eGFR_fu,10)

```

	<code>id</code>	<code>avg_eGFR</code>	<code>length_of_followup</code>
1:	1	43.04333	53.8341
2:	2	38.93294	15.7618
3:	3	85.72000	34.7487
4:	4	76.59308	31.8055
5:	5	NA	192.1423
6:	6	85.66435	89.7441
7:	7	64.21758	131.7151
8:	8	66.28333	9.1226
9:	9	86.35750	26.5627
10:	10	107.00429	20.4599

```

1 # Now let us specify the given range labels that will be used later on.
2 ranges<-c('(0,15]', '(15,30]', '(30,60]', '(60,90]', '(90,max(egfr))')
3
4 # Now let us count and store the number of patients in each of those ranges.
5 num_patients<-c(nrow(Avg_eGFR_fu[avg_eGFR>0 & avg_eGFR<=15]),

```

```

6      nrow(Avg_eGFR_fu[avg_eGFR>15 & avg_eGFR<=30]),
7      nrow(Avg_eGFR_fu[avg_eGFR>30 & avg_eGFR<=60]),
8      nrow(Avg_eGFR_fu[avg_eGFR>60 & avg_eGFR<=90]),
9      nrow(Avg_eGFR_fu[avg_eGFR>90]))
10
11 # Now let us form a data table to better visualize the number of rows or
12 # samples in each range using the 2 vectors we computed earlier.
13 egfr_ranges<-data.table(ranges,num_patients)
14 egfr_ranges

          ranges num_patients
1:      (0,15]           1
2:      (15,30]          9
3:      (30,60]         83
4:      (60,90]         82
5: (90,max(egfr))       36

1 # Now let us print the number of patients with missing eGFR using cat().
2 # This can simply be done by counting the number of rows with avg_eGFR as
3 # NA using nrow() and using is.na() to find the NA rows.
4 cat('The number of patients with missing average eGFR is',
5     nrow(Avg_eGFR_fu[is.na(avg_eGFR)]))

```

The number of patients with missing average eGFR is 39

Problem 1.c (6 points)

- For patients with average eGFR in the (90,max(eGFR)) range, collect their identifier, sex, age at baseline, average eGFR, time of last eGFR reading and number of eGFR measurements taken in a data table.
- Print the summary of the new dataset.

```

1 # The next objective is to collect specific data of the given patients.
2 # For convenience, let us add the average eGFR that we computed earlier
3 # for Avg_eGFR_fu to the original data table lgfr121.
4 lgfr121<-lgfr12[, `:=` (avg_eGFR=mean(egfr)),
5                     by='id']
6

```

```

7 # Now we use lgfr121 to return a set of specific details of the patients
8 # whose ids are greater than 90. For a given id, sex, baseline age and average
9 # eGFR are all constant. So we simply use the unique() function to return
10 # a single value for each.
11
12 # The number of measurements by id can simply be returned using .N operator.
13 # Finally, the last eGFR reading by id can be returned using .SD and
14 # specifying index of row as .N. This retrieves the last row by id which
15 # had required fu.years value.
16 lgfr122<-lgfr121[avg_eGFR>90,.(sex=unique(sex),
17                                baseline_age=unique(baseline.age),
18                                avg_eGFR=unique(avg_eGFR),
19                                last_reading_time=.SD[.N,fu.years],
20                                Nmeasurements=.N),by='id']
21
22 # Let us print the summary of the new data table lgfr122 using summary().
23 summary(lgfr122)

```

	id	sex	baseline_age	avg_eGFR
Min.	: 10.00	Min. :0.0000	Min. :22.10	Min. : 90.04
1st Qu.	: 79.75	1st Qu.:0.0000	1st Qu.:48.05	1st Qu.: 94.49
Median	:137.00	Median :0.0000	Median :56.45	Median :105.61
Mean	:139.19	Mean :0.3889	Mean :57.54	Mean :105.65
3rd Qu.	:215.75	3rd Qu.:1.0000	3rd Qu.:67.83	3rd Qu.:113.91
Max.	:250.00	Max. :1.0000	Max. :90.90	Max. :147.69
last_reading_time		Nmeasurements		
Min.	:0.000	Min. : 1.00		
1st Qu.	:1.275	1st Qu.: 4.00		
Median	:3.469	Median : 7.00		
Mean	:3.253	Mean :11.11		
3rd Qu.	:5.201	3rd Qu.:10.00		
Max.	:6.590	Max. :57.00		

```

1 # Finally, we print the first 10 rows of the data table using head().
2 head(lgfr122,10)

```

	id	sex	baseline_age	avg_eGFR	last_reading_time	Nmeasurements
1:	10	0	50.4	107.00429	5.1964	7
2:	14	0	65.1	116.09200	4.0986	10
3:	25	0	40.1	95.35625	4.2847	8

4:	31	0	74.8	113.59250	1.4675	8
5:	33	0	74.2	116.35000	1.6016	4
6:	45	1	24.9	91.25000	0.0000	1
7:	49	1	68.2	128.25800	6.1602	5
8:	52	1	56.3	93.31544	6.4805	57
9:	79	0	65.6	91.45057	5.2156	35
10:	80	0	67.7	106.09600	2.2834	5

Problem 1.d (9 points)

For patients 3, 37, 162 and 223:

- Plot the patient's eGFR measurements as a function of time.
- Fit a linear regression model and add the regression line to the plot.
- Report the 95% confidence interval for the regression coefficients of the fitted model.
- Using a different colour, plot a second regression line computed after removing the extreme eGFR values (one each of the highest and the lowest value).

(All plots should be displayed in the same figure. The plots should be appropriately labelled and the results should be accompanied by some explanation as you would communicate it to a colleague with a medical background with a very little statistical knowledge.)

```

1  # Next we will be plotting egfr v time. We will also be plotting two regression lines
2  # for the data points for the given ids. The first with all data points and the
3  # second with removing the extreme values. The 1st line is in black while
4  # the 2nd is in purple.
5
6  # First let's partition the screen into 4 parts using par().
7  par(mfrow=c(2,2))
8  # The vector of ids to plot.
9  a<-c(3,37,162,223)
10
11 for(i in a){
12   # Next for each id:
13   # we plot egfr as a function of time.
14   p1<-plot(lgfr12[id==i]$fu.years,lgfr12[id==i]$egfr,
15             xlab = 'time',ylab = 'egfr',
16             main = paste('eGFR measurements as a \nfunction of time (id =',i,')'))
17   # Now we model the first linear regression line for egfr and fu.years using lm().
18   regr1<-lm(egfr ~ fu.years,data = lgfr12[id==i])
19   # We then pass the regression object regr1 to abline() that draws the

```

```

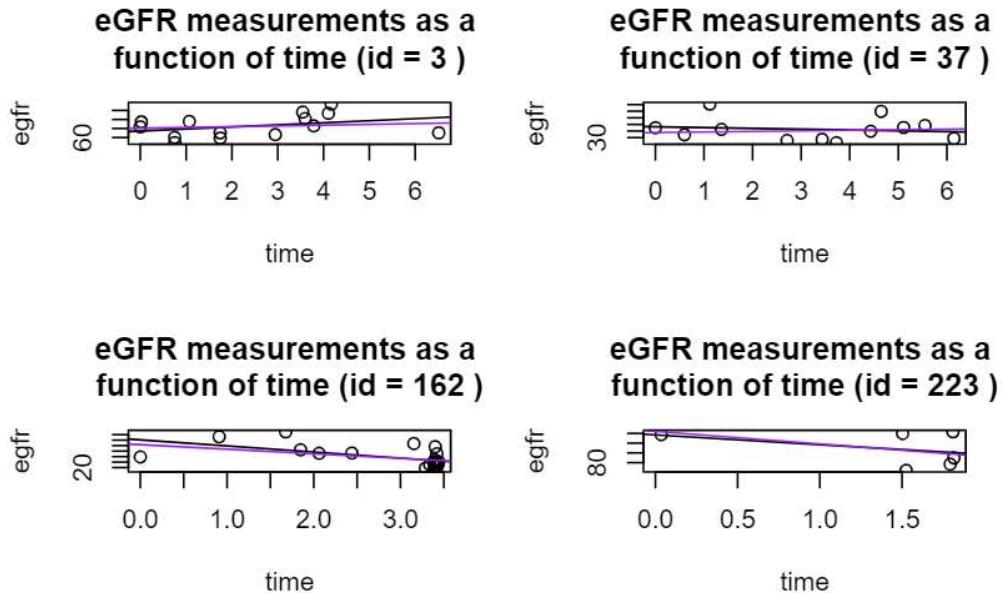
20  # regression line. This line is added to the previous scatterplot.
21  p1<-p1+abline(regr1)
22  # Now we report 95% confidence value using confint() function.
23  cat('\n95% confidence interval (id = ',i,')\n',
24      confint(regr1,'fu.years',level = 0.95),'\n')
25  # We now find out the extreme values i.e, the maximum and minimum
26  # values using max() and min() respectively. To exclude NA values
27  # from the computation, we set na.rm=TRUE.
28  max_egfr<-max(lgfr12[id==i]$egfr,na.rm = TRUE)
29  min_egfr<-min(lgfr12[id==i]$egfr,na.rm = TRUE)
30  # Now filter rows by id
31  new_lgfr3<-lgfr12[id==i]
32  # Next exclude the extreme values from the filtered rows.
33  new_lgfr3<-new_lgfr3[egfr!=max_egfr & egfr!=min_egfr]
34  # Now create a second linear regression model with
35  # the new values.
36  regr11<-lm(egfr ~ fu.years,data = new_lgfr3)
37  # The resultant regression line is added to the final plot.
38  p1<-p1+abline(regr11,col='purple')
39  # Print resultant plot (Repeat process for each id)
40  p1
41 }

```

95% confidence interval (id = 3)
-3.151128 12.25612

95% confidence interval (id = 37)
-3.595705 2.37859

95% confidence interval (id = 162)
-9.257727 -1.872262



95% confidence interval (id = 223)

-85.93757 45.9659

```

1  #### Explanation of the plots:
2  # The regression plots for all 4 ids can be visualized as shown below.
3  # The black line indicates the first regression line on all data points
4  # by given id. The purple line indicates the second regression line.
5  # This line was formed out of the same data points as before except we
6  # exclude extreme values. The ids are 3,37,162 and 223 respectively.
7  # Also observe how the slopes of the new regression lines (in purple)
8  # changed after we removed extreme values (or outliers).
9
10 # On removing these values, one can expect that the correlation between
11 # the variables would increase, be it much more positive (tending to 1)
12 # or much more negative (tending to -1).
13
14 # We also observe a change in slope of the line. As the model does not
15 # try to fit-in those extreme values anymore, we observe a increase or
16 # decrease in slope accordingly. This can be observed in the plots below.

```

Problem 2 (25 points)

The MDRD4 and CKD-EPI equations are two different ways of estimating the glomerular filtration rate (eGFR) in adults:

$$\text{MDRD4} = 175 \times (\text{SCR})^{-1.154} \times \text{AGE}^{-0.203} [\times 0.742 \text{ if female}] [\times 1.212 \text{ if black}]$$

, and

$$\text{CKD-EPI} = 141 \times \min(\text{SCR}/\kappa, 1)^\alpha \times \max(\text{SCR}/\kappa, 1)^{-1.209} \times 0.993^{\text{AGE}} [\times 1.018 \text{ if female}] [\times 1.159 \text{ if black}]$$

, where:

- SCR is serum creatinine (in mg/dL)
- κ is 0.7 for females and 0.9 for males
- α is -0.329 for females and -0.411 for males

Problem 2.a (7 points)

For the `scr.csv` dataset,

- Examine a summary of the distribution of serum creatinine and report the inter-quartile range.
- If you suspect that some serum creatinine values may have been reported in $\mu\text{mol/L}$ convert them to mg/dL by dividing by 88.42.
- Justify your choice of values to convert and examine the distribution of serum creatinine following any changes you have made.

```

1 # First lets read the required file using fread().
2 scr_d<-fread('scr.csv')
3 # Then convert it to a data table using setDT().
4 scr_d<-setDT(scr_d)
5 # Now order the data table by scr values using the with() function.
6 scr_d<-scr_d[with(scr_d,order(scr))]
7 # Next, print a summary of the serum creatinine values:
8 summary(scr_d$scr)

```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
	0.400	0.900	1.300	3.072	2.800	76.000	18

```

1 # Delete NA values using na.omit()
2 scr_d<-na.omit(scr_d)
3 # Now lets store the first and third quartile values in q.
4 # This will be useful later on.
5 q<-quantile(scr_d$scr,c(0.25,0.75))
6 # The inter-quartile range is computed using the IQR() function.
7 iqr<-IQR(scr_d$scr)
8 # Print the IQR using cat().
9 cat('The Inter-Quartile Range is',iqr)

```

The Inter-Quartile Range is 1.9

```

1 # Now the upper and lower bound for the whiskers can be computed as follows:
2 # LB = Q1 - 1.5*IQR and UB = Q3 + 1.5*IQR
3 lb<-q[1]-1.5*iqr
4 ub<-q[2]+1.5*iqr
5 # Now the assumption is that the values above ub or below lb are outliers.
6 # In this case, they are the values to be scaled to mg/dL.
7 scr_d_t<-scr_d[, 
8   scr:=ifelse(scr>ub,
9             scr/88.42,
10            scr)]
11 #### Choice of Values:
12 # There are some scr values that may have been reported in umol/litre.
13 # One easy way to find out which ones is to simply find the "outliers"
14 # in a boxplot. The data is first sorted in ascending order and then compute
15 # the lower bound and upper bound using the following formulae:
16 # LB = Q1 - 1.5*IQR and UB = Q3 + 1.5*IQR
17 # From here on, pretty much any value lesser than LB or greater than UB
18 # is assumed to be an outlier and thus the values that are in umol/litre
19 # and are thus converted to mg/dL by dividing by 88.42.

```

Problem 2.b (11 points)

- Compute the eGFR according to the two equations using the newly converted SCR values.
- Report (rounded to the second decimal place) mean and standard deviation of the two eGFR vectors and their Pearson correlation coefficient.
- Report the same quantities according to strata of MDRD4 eGFR: (0 – 60), (60 – 90) and (> 90).

- Print first 15 values for both datasets using `head()`.

```

1 # The two functions below compute the new scr values according to the 2
2 # equations specified earlier.
3 # MDRD4_calculate() calculates new scr values by equation 1.
4 MDRD4_calculate<-function(scr,sex,ethnic,age){
5   e_factor<-ifelse(ethnic=='Black',1.212,1)
6   f_factor<-ifelse(sex=='Female',0.742,1)
7   eq<-175*scr^(-1.154)*age^(-0.203)*f_factor*e_factor
8   return(eq)
9 }
10 # Call the function on our data table scr_d_t and pass required parameters.
11 mdrd4<-MDRD4_calculate(scr_d_t$scr,scr_d_t$sex,scr_d_t$ethnic,scr_d_t$age)
12 # Now report the mean and standard deviation of the new scr values.
13 # While the mean is the 4th value in the summary, the Standard deviation is
14 # computed using sd(). The results are rounded to 2 decimal places using round().
15 cat('Mean of mdrd4 scr values:',round(summary(mdrd4)[4],2),'\n')

```

Mean of mdrd4 scr values: 188.98

```
1 cat('Standard Deviation of mdrd4 scr values:',round(sd(mdrd4),2),'\n\n') #na
```

Standard Deviation of mdrd4 scr values: 359.03

```

1 # CKD_EPI_calculate() calculates new scr values by equation 2.
2 CKD_EPI_calculate<-function(scr,sex,ethnic,age){
3
4   l<-length(scr)
5   K<-ifelse(sex=='Female',0.7,0.9)
6   A<-ifelse(sex=='Female',-0.329,-0.411)
7   e_factor<-ifelse(ethnic=='Black',1.159,1)
8   f_factor<-ifelse(sex=='Female',1.018,1)
9
10  min_factor<-max_factor<-rep(0,1)
11
12  for(i in 1:l){
13    min_factor[i]<-min(scr[i]/K[i],1)
14    max_factor[i]<-max(scr[i]/K[i],1)
15  }

```

```

16   eq<-141*min_factor^(A)*max_factor^(-1.209)*0.993^(age)*f_factor*e_factor
17   return(eq)
18 }
19 # Call the function on our data table scr_d_t and pass required parameters.
20 ckd_epi<-CKD_EPI_calculate(scr_d_t$scr,scr_d_t$sex,scr_d_t$ethnic,scr_d_t$age)
21 # Report mean and standard deviation like before and also the Pearson
22 # Correlation coefficient using cor() function.
23 cat('Mean of ckd epi scr values:',round(summary(ckd_epi)[4],2),'\n')

```

Mean of ckd epi scr values: 85.84

```
1 cat('Standard Deviation of ckd epi scr values:',round(sd(ckd_epi),2),'\n')
```

Standard Deviation of ckd epi scr values: 64.27

```
1 cat('Pearson correlation coefficient:',round(cor(mdrd4,ckd_epi),2),'\n\n')
```

Pearson correlation coefficient: 0.86

```

1 # Form a new data table new_scr_values using the values computed earlier.
2 new_scr_values<-data.table('MDRD4'=mdrd4,'CKD-EPI'=ckd_epi)
3 # Now we filter the values by several strata of MDRD4.
4 s1<-new_scr_values[MDRD4>0 & MDRD4<60]
5 s2<-new_scr_values[MDRD4>60 & MDRD4<90]
6 s3<-new_scr_values[MDRD4>90]
7
8 # Now we print the mean, standard deviation and correlation coefficient for each strata.
9 cat('MDRD4 (0-60 Strata):\nMDRD4 Equation Mean:'
10   ,round(summary(s1$MDRD4)[4],2),
11   '\nCKD-EPI Equation Mean:'
12   ,round(summary(s1$`CKD-EPI`)[4],2),
13   '\nPearson Coefficient:'
14   ,round(cor(s1$MDRD4,s1$`CKD-EPI`),2),'\n\n')

```

MDRD4 (0-60 Strata):
 MDRD4 Equation Mean: 31.9
 CKD-EPI Equation Mean: 33.15
 Pearson Coefficient: 0.99

```
1 cat('MDRD4 (60-90 Strata):\nMDRD4 Equation Mean:\n',
2     ,round(summary(s2$MDRD4)[4],2),
3   '\nCKD-EPI Equation Mean:\n',
4   ,round(summary(s2$`CKD-EPI`)[4],2),
5   '\nPearson Coefficient:\n',
6   ,round(cor(s2$MDRD4,s2$`CKD-EPI`),2),'\n\n')
```

MDRD4 (60-90 Strata):
 MDRD4 Equation Mean: 73.41
 CKD-EPI Equation Mean: 80.18
 Pearson Coefficient: 0.93

```
1 cat('MDRD4 (>90 Strata):\nMean:\n',
2     ,round(summary(s3$MDRD4)[4],2),
3   '\nCKD-EPI Equation Mean:\n',
4   ,round(summary(s3$`CKD-EPI`)[4],2),
5   '\nPearson Coefficient:\n',
6   ,round(cor(s3$MDRD4,s3$`CKD-EPI`),2),'\n\n')
```

MDRD4 (>90 Strata):
 Mean: 466.01
 CKD-EPI Equation Mean: 156.02
 Pearson Coefficient: 0.95

```
1 # Print first 15 values of new data table.
2 head(new_scr_values,15)
```

MDRD4	CKD-EPI
1: 184.9777	137.81686
2: 167.9425	115.32906
3: 169.6141	117.78528
4: 135.3041	119.37447
5: 195.2410	145.41646

```

6: 144.8688 129.87348
7: 120.9417 96.01479
8: 132.2497 115.25445
9: 170.1717 145.32828
10: 147.8353 113.65115
11: 169.6141 117.78528
12: 150.3311 134.51609
13: 202.6025 150.61469
14: 134.0324 117.70908
15: 138.0812 122.77628

```

Problem 2.c (7 points)

- Produce a scatter plot of the two eGFR vectors, and add vertical and horizontal lines (i.e.) corresponding to median, first and third quantiles.
- Is the relationship between the two eGFR equations linear? Justify your answer.

```

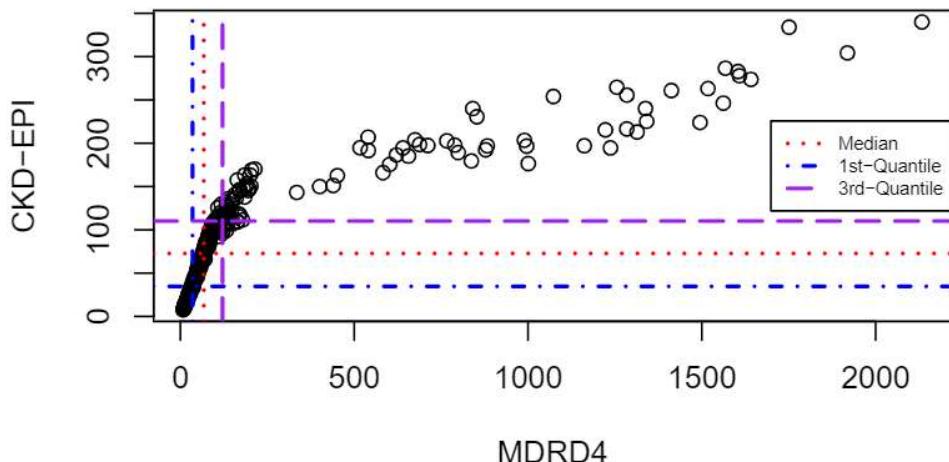
1 # Let's first retrieve the median, 1st and 3rd quartile from the summary
2 # of the 2 new scr values
3 # (They are in the 3rd, 2nd and 5th indicies respectively):
4 mdrd4_v<-summary(new_scr_values$MDRD4)[c(3,2,5)]
5 ckd_v<-summary(new_scr_values$`CKD-EPI`)[c(3,2,5)]
6 # Next, a scatter plot is plotted between the 2 set of values. The title,
7 # xlabel, ylabel can be specified using additional parameters as follows:
8 plot(x=new_scr_values$MDRD4,y=new_scr_values$`CKD-EPI`,
9       xlab = 'MDRD4',ylab='CKD-EPI',
10      main = 'Scatter Plot between the two new SCR equations')
11 # Now we add three lines.
12 # These indicate median, 1st and 3rd quartile respectively.
13 # We do the same for MDRD4 and then for CKD-EPI values that give us 6
14 # lines overall (3 Vertical and 3 Horizontal).
15 # The width of each line and type of line can be customized
16 # using lty and lwd.
17 abline(v=mdrd4_v,col=c('red','blue','purple'),
18         lty=c(3,4,5),lwd=c(2,2,2))
19 abline(h=ckd_v,col=c('red','blue','purple'),
20         lty=c(3,4,5),lwd=c(2,2,2))
21 # Then a legend is added using legend() to indicate the values
22 # each line points to.
23 legend('right',legend=c('Median','1st-Quantile','3rd-Quantile'),
24        col=c('red','blue','purple')),

```

25

```
lty=c(3,4,5),lwd=c(2,2,2),cex=0.6)
```

Scatter Plot between the two new SCR equations



```
1  ### The relationship between the equations:  
2  # The relationship between the two eGFR equations appears to follow a  
3  # linear trend until we go past the third quantile. From there,  
4  # we clearly observe that the relationship is not linear anymore.  
5  # We see the points scattered with increase in MDRD4 and overall  
6  # it appears as if it is forming a curve thus hinting at a quadratic  
7  # relationship.
```

Problem 3 (31 points)

You have been provided with electronic health record data from a study cohort. Three CSV (Comma Separated Variable) files are provided on learn.

The first file is a cohort description file `cohort.csv` file with fields:

- `id` = study identifier
- `yob` = year of birth
- `age` = age at measurement
- `bp` = systolic blood pressure
- `albumin` = last known albuminuric status (categorical)
- `diabetes` = diabetes status

The second file `lab1.csv` is provided by a laboratory after measuring various biochemistry levels in the cohort blood samples. Notice that a separate lab identifier is used to anonymise results from the cohort. The year of birth is also provided as a check that the year of birth aligns between the two merged sets.

- `LABID` = lab identifier
- `yob` = year of birth
- `urea` = blood urea
- `creatinine` = serum creatinine
- `glucose` = random blood glucose

To link the two data files together, a third linker file `linker.csv` is provided. The linker file includes a `LABID` identifier and the corresponding cohort `id` for each person in the cohort.

Problem 3.a (6 points)

- Using all three files provided on learn, load and merge to create a single data table based dataset `cohort.dt`. This will be used in your analysis.
- Perform assertion checks to ensure that all identifiers in `cohort.csv` have been accounted for in the final table and that any validation fields are consistent between sets.
- After the checks are complete, drop the identifier that originated from `lab1.csv` dataset `LABID`.
- Ensure that a single `yob` field remains and rename it to `yob`.
- Ensure that the `albumin` field is converted to a factor and the ordering of the factor is `1="normo", 2="micro", 3="macro"`.
- Print first 10 values of the new dataset using `head()`.

```

1 # Read the 3 datasets.
2 d1<-fread('cohort.csv')
3 d2<-fread('lab1.csv')
4 d3<-fread('linker.csv')
5 # Convert datasets to data tables.
6 d1<-setDT(d1)
7 d2<-setDT(d2)
8 d3<-setDT(d3)
9 # Next, we merge the data tables by key column.
10 # First d1 and d3 are merged by common key column 'id'.
11 d4<-merge.data.table(d1,d3,by.x = 'id',by.y = 'id',all = TRUE)
12 # Next we merge the resultant data table d4 with d2
13 # by c('LABID', 'yob') to create final data table cohort.dt.
14 cohort.dt<-merge.data.table(d4,d2,
15                               by.x = c('LABID','yob'),
16                               by.y = c('LABID','yob'),all = TRUE)

1 # Next we perform a simple assertion check to check whether all ids
2 # in d1 were accounted for in the final data table.
3 # The idea is that if the intersection of d1 and cohort.dt ids gives
4 # number of ids in d1, that simply means that no ids were lost in
5 # merge process.
6 length(intersect(d1$id,cohort.dt$id))==length(d1$id) # returns TRUE

```

[1] TRUE

```

1 # Next we delete a row called LABID by initializing it to NULL.
2 cohort.dt[,LABID:=NULL]
3
4 # The categorical variable albumin is converted to numeric by following
5 # normo=1,micro=2,macro=3 convention.
6 cohort.dt$albumin<-c(normo=1,micro=2,macro=3)[cohort.dt$albumin]
7 # Print first 10 rows.
8 head(cohort.dt,10)

```

	yob	id	age	bp	diabetes	albumin	urea	creatinine	glucose	
1:	1986	PID_285	33	80		0	1	37.0	106.104	100
2:	1980	PID_153	39	70		1	1	20.0	70.736	121
3:	1951	PID_13	68	70		1	2	72.0	185.682	208

4:	1965	PID_110	54	70	1	NA	50.1	167.998	233
5:	1953	PID_222	66	70	1	2	30.0	150.314	248
6:	2002	PID_103	17	60	0	1	32.0	185.682	92
7:	1954	PID_200	65	80	0	1	37.0	132.630	92
8:	1955	PID_378	64	70	0	1	27.0	61.894	97
9:	1964	PID_267	55	80	0	1	17.0	106.104	133
10:	1996	PID_271	23	80	0	1	34.0	97.262	111

Problem 3.b (10 points)

- Create a copy of the dataset where you will impute all missing values.
- Update any missing age fields using the year of birth.
- Perform mean imputation for all other continuous variables by writing a single function called `impute.to.mean()` and impute to mean, impute any categorical variable to the mode.
- Print first 15 values of the new dataset using `head()`.
- Compare each distribution of the imputed and non-imputed variables and decide which ones to keep for further analysis. Justify your answer.

```

1 # Now let us create a copy of the dataset using cohort.dt %>% copy().
2 # We also fill missing values by calculating age using corresponding
3 # yob as follows:
4 cohort.dt.impute<-cohort.dt %>% copy() %>%
5   .[, age:=ifelse(is.na(age), 2019-yob,age)]
6 # A function is created which fills the missing values of a
7 # column with the mean.
8 impute.to.mean <- function(x) {
9   # find which values are missing
10  na.idx <- is.na(x)
11  # replace NAs with the mean computed over observed values
12  x[na.idx] <- mean(x, na.rm = TRUE)
13  # return the vector with imputed mean values
14  return(x)
15 }
16 # A function to calculate mode of a column.
17 mode_exclude_na <- function(x) {
18   uq <- na.omit(unique(x))
19   t <- tabulate(match(x, uq))
20   uq[t == max(t)]
21 }
22 # A function is created which fills the missing values

```

```

23 # of a column with the mode.
24 impute.to.mode <- function(x) {
25     # find which values are missing
26     na.idx <- is.na(x)
27     # replace NAs with the mean computed over observed values
28     x[na.idx] <- mode_exclude_na(x)
29     # return the vector with imputed values
30     return(x)
31 }
32 # Impute missing values in bp, urea, creatinine and glucose
33 # columns with mean
34 cohort.dt.impute$bp <- impute.to.mean(cohort.dt.impute$bp)
35 cohort.dt.impute$urea <- impute.to.mean(cohort.dt.impute$urea)
36 cohort.dt.impute$creatinine<- impute.to.mean(cohort.dt.impute$creatinine)
37 cohort.dt.impute$glucose<- impute.to.mean(cohort.dt.impute$glucose)
38 # Impute missing values in albumin and diabetes columns with mode
39 cohort.dt.impute$albumin<- impute.to.mode(cohort.dt.impute$albumin)
40 cohort.dt.impute$diabetes<- impute.to.mode(cohort.dt.impute$diabetes)
41 # Return first 15 values.
42 head(cohort.dt.impute,15)

```

	yob	id	age	bp	diabetes	albumin	urea	creatinine	glucose
1:	1986	PID_285	33	80	0	1	37.00000	106.1040	100.0000
2:	1980	PID_153	39	70	1	1	20.00000	70.7360	121.0000
3:	1951	PID_13	68	70	1	2	72.00000	185.6820	208.0000
4:	1965	PID_110	54	70	1	1	50.10000	167.9980	233.0000
5:	1953	PID_222	66	70	1	2	30.00000	150.3140	248.0000
6:	2002	PID_103	17	60	0	1	32.00000	185.6820	92.0000
7:	1954	PID_200	65	80	0	1	37.00000	132.6300	92.0000
8:	1955	PID_378	64	70	0	1	27.00000	61.8940	97.0000
9:	1964	PID_267	55	80	0	1	17.00000	106.1040	133.0000
10:	1996	PID_271	23	80	0	1	34.00000	97.2620	111.0000
11:	1964	PID_105	55	90	1	1	88.00000	176.8400	143.0000
12:	1940	PID_375	79	80	0	1	44.00000	106.1040	111.0000
13:	1961	PID_89	58	110	0	3	52.00000	194.5240	251.0000
14:	1998	PID_24	21	70	0	1	57.42572	271.6664	148.0365
15:	1981	PID_349	38	80	0	1	19.00000	44.2100	99.0000

1 #Answer in this chunk

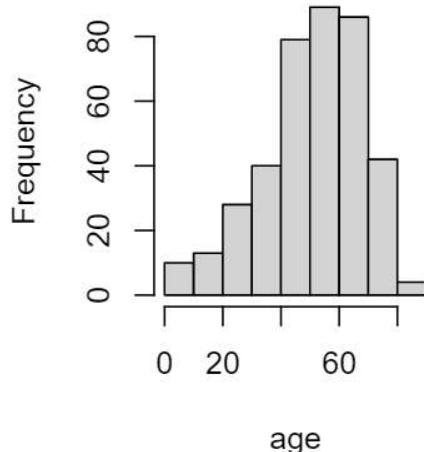
```

1  # Store the imputed values in a data table dx2.
2  dx2<-data.table(age=cohort.dt.impute$age,
3                    bp=cohort.dt.impute$bp,
4                    urea=cohort.dt.impute$urea,
5                    creatinine =cohort.dt.impute$creatinine,
6                    albumin=cohort.dt.impute$albumin,
7                    glucose=cohort.dt.impute$glucose,
8                    diabetes=cohort.dt.impute$diabetes)
9  # Store the corresponding initial or Non-imputed values in a data table dx1.
10 dx1<-data.table(age=cohort.dt$age,
11                   bp=cohort.dt$bp,
12                   urea=cohort.dt$urea,
13                   creatinine =cohort.dt$creatinine,
14                   albumin=cohort.dt$albumin,
15                   glucose=cohort.dt$glucose,
16                   diabetes=cohort.dt.impute$diabetes)
17
18 # Function to plot distributions of imputed variables.
19 sbsplots_i <- function(varname, vars){
20   # Partition screen into 2 parts
21   par(mfrow = c(1,2))
22   # Plot histogram using hist()
23   hist(vars[,varname],
24         main = paste0("Histogram of \nImputed ", varname),
25         xlab = varname)
26   # Plot boxplot using boxplot()
27   boxplot(vars[,varname],
28            main = paste0("Boxplot of \nImputed ", varname),
29            xlab = varname,
30            ylab = "Value")
31 }
32 # Function to plot distributions of imputed variables.
33 sbsplots_ni <- function(varname, vars){
34   # Partition screen into 2 parts
35   par(mfrow = c(1,2))
36   # Plot histogram using hist()
37   hist(vars[,varname],
38         main = paste0("Histogram of \nNon-Imputed ", varname),
39         xlab = varname)
40   # Plot boxplot using boxplot()
41   boxplot(vars[,varname],

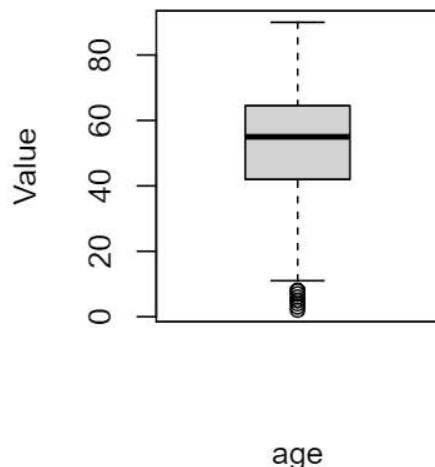
```

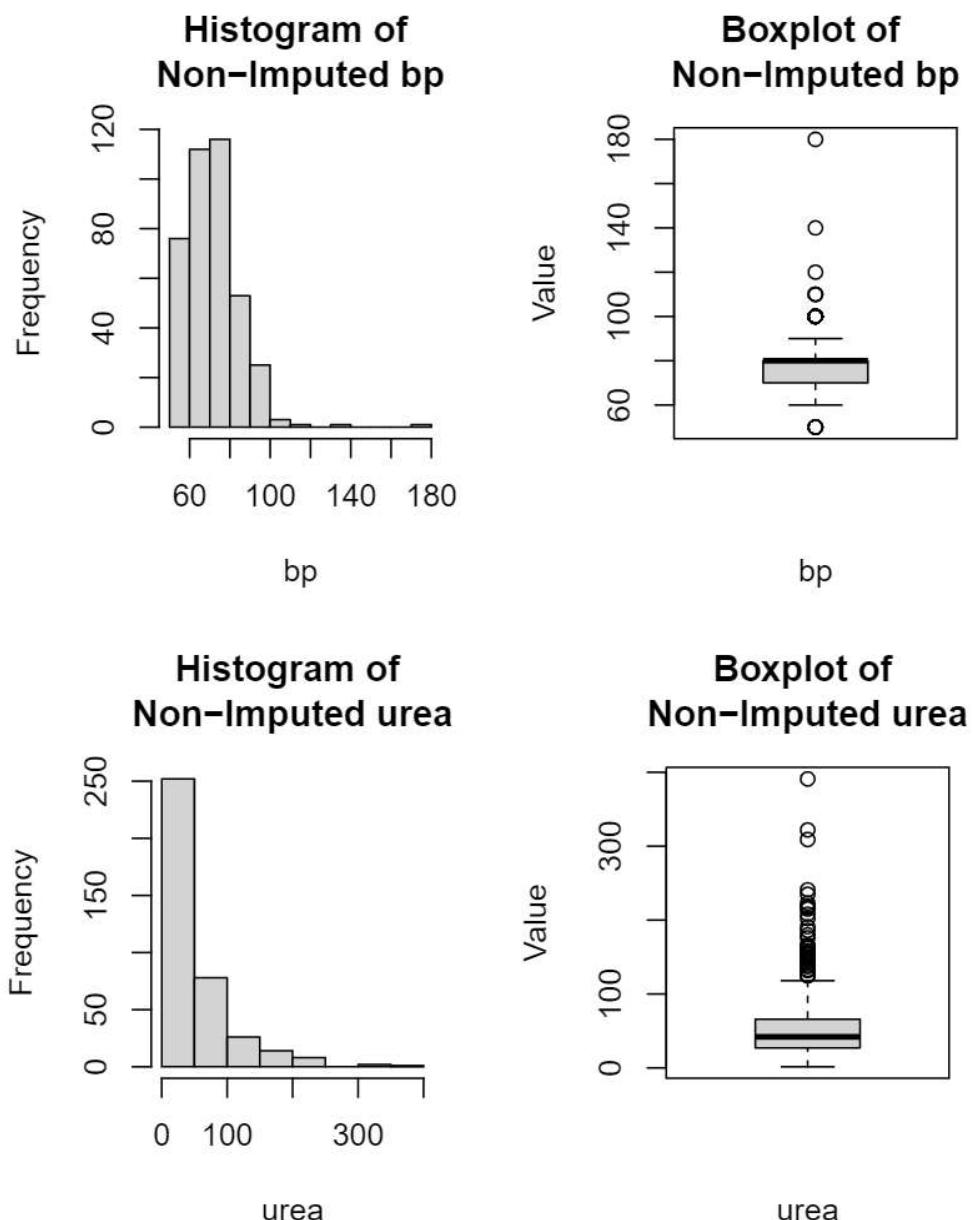
```
42     main = paste0("Boxplot of \nNon-Imputed ", varname),  
43     xlab = varname,  
44     ylab = "Value")  
45 }  
  
1 # Store required Non-imputed column names  
2 numcols1 <- dx1 %>% colnames  
3 # Use sapply() to apply sbsplots_ni() to required columns.  
4 sapply(numcols1, sbsplots_ni, vars = data.frame(dx1)) # Non-Imputed
```

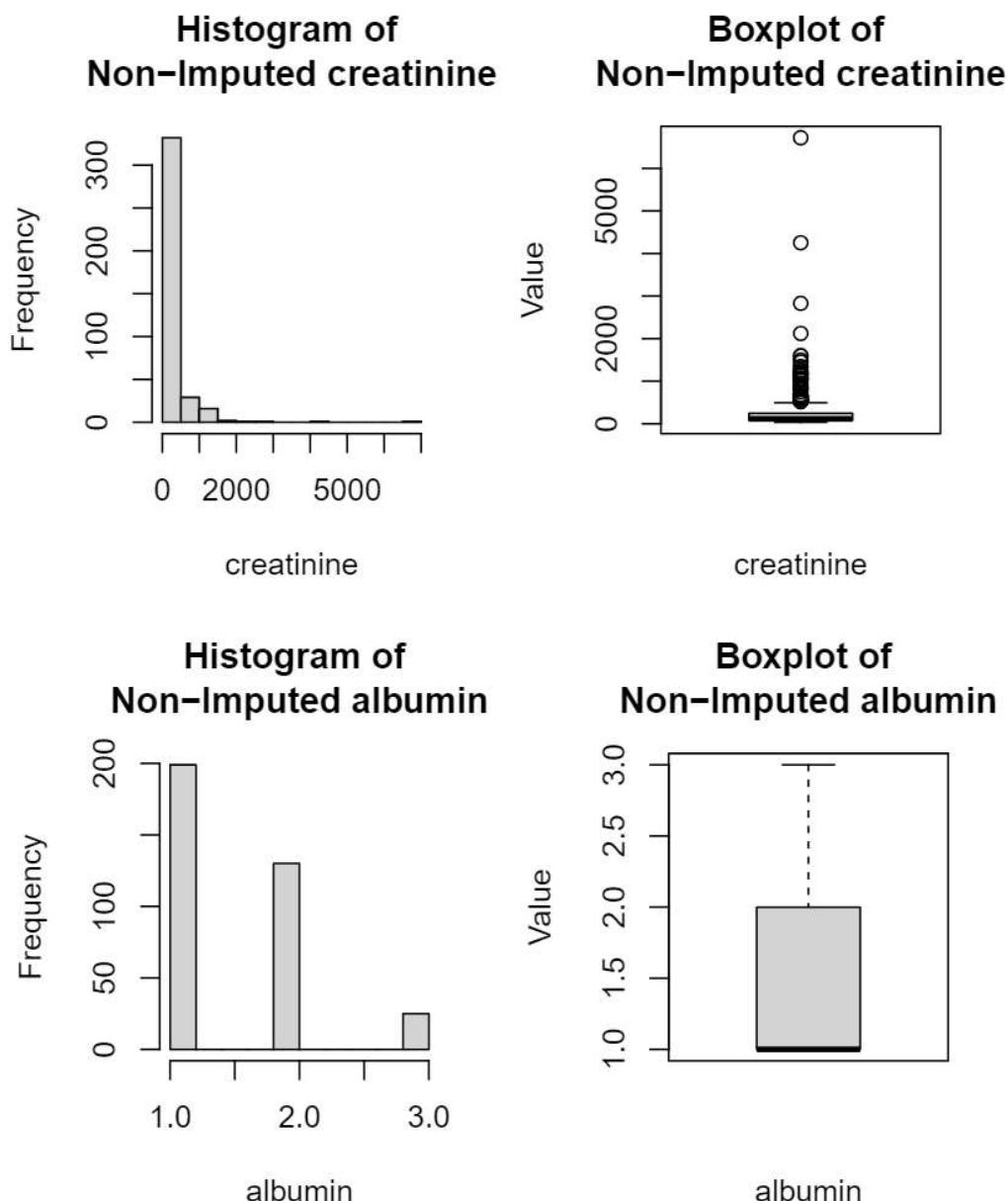
Histogram of Non-Imputed age

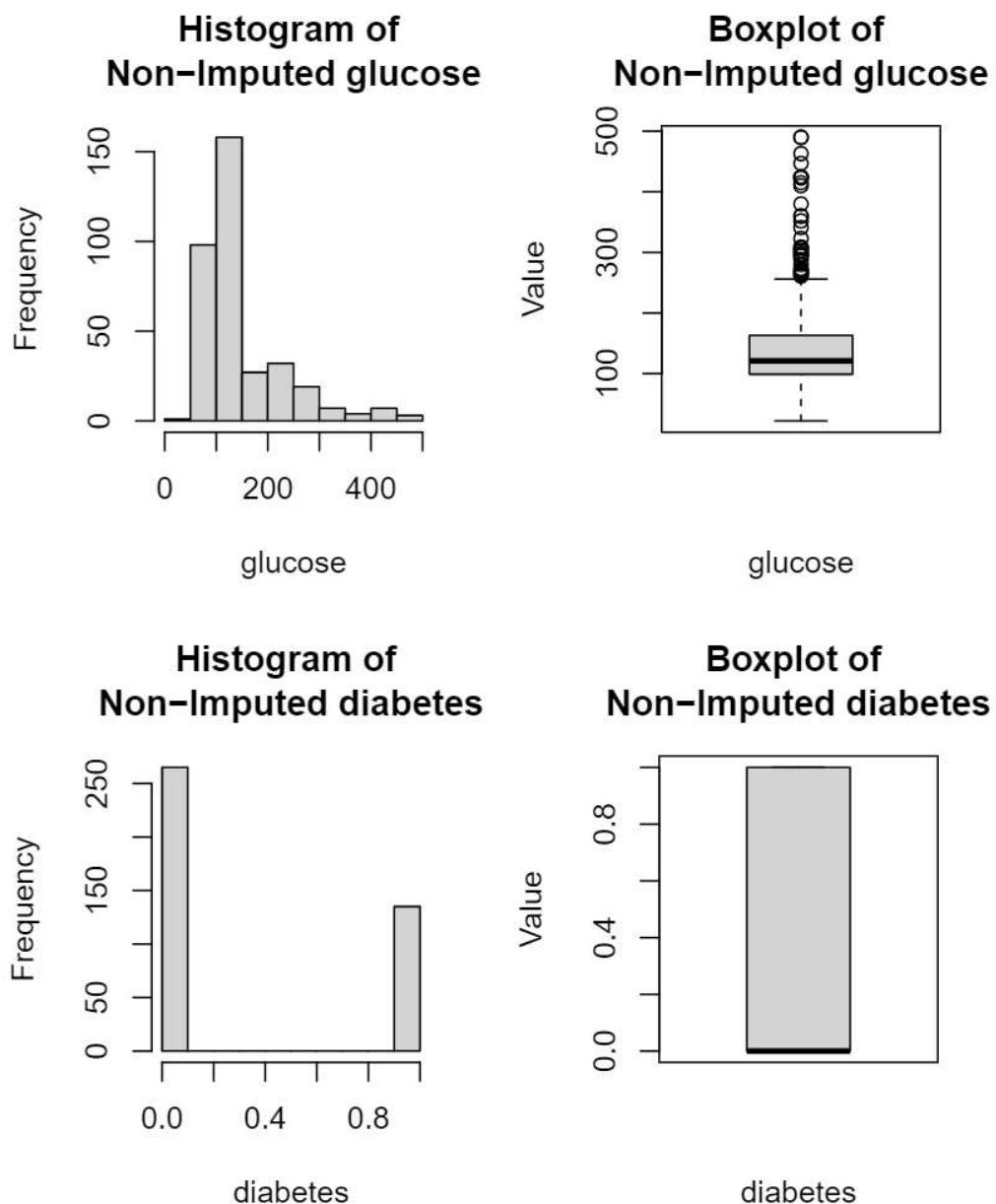


Boxplot of Non-Imputed age







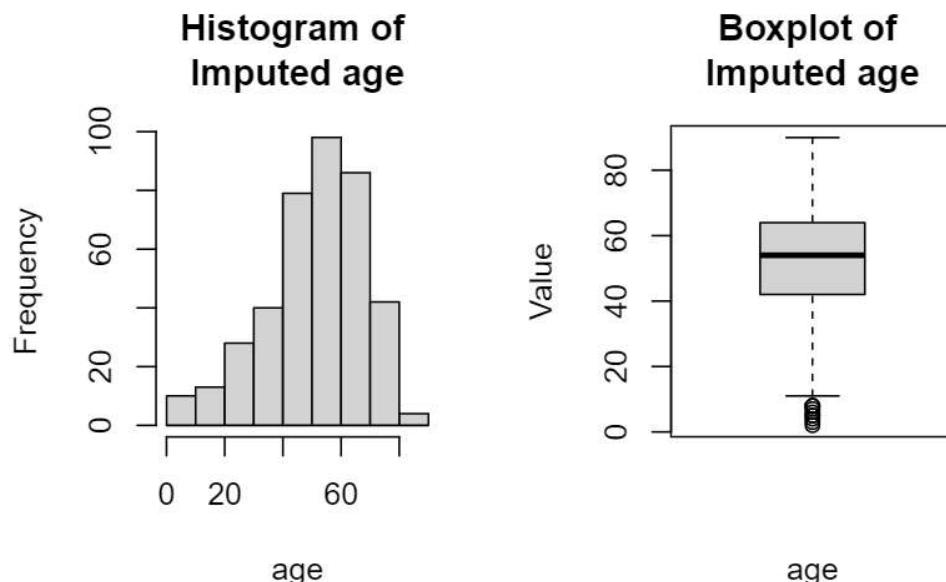


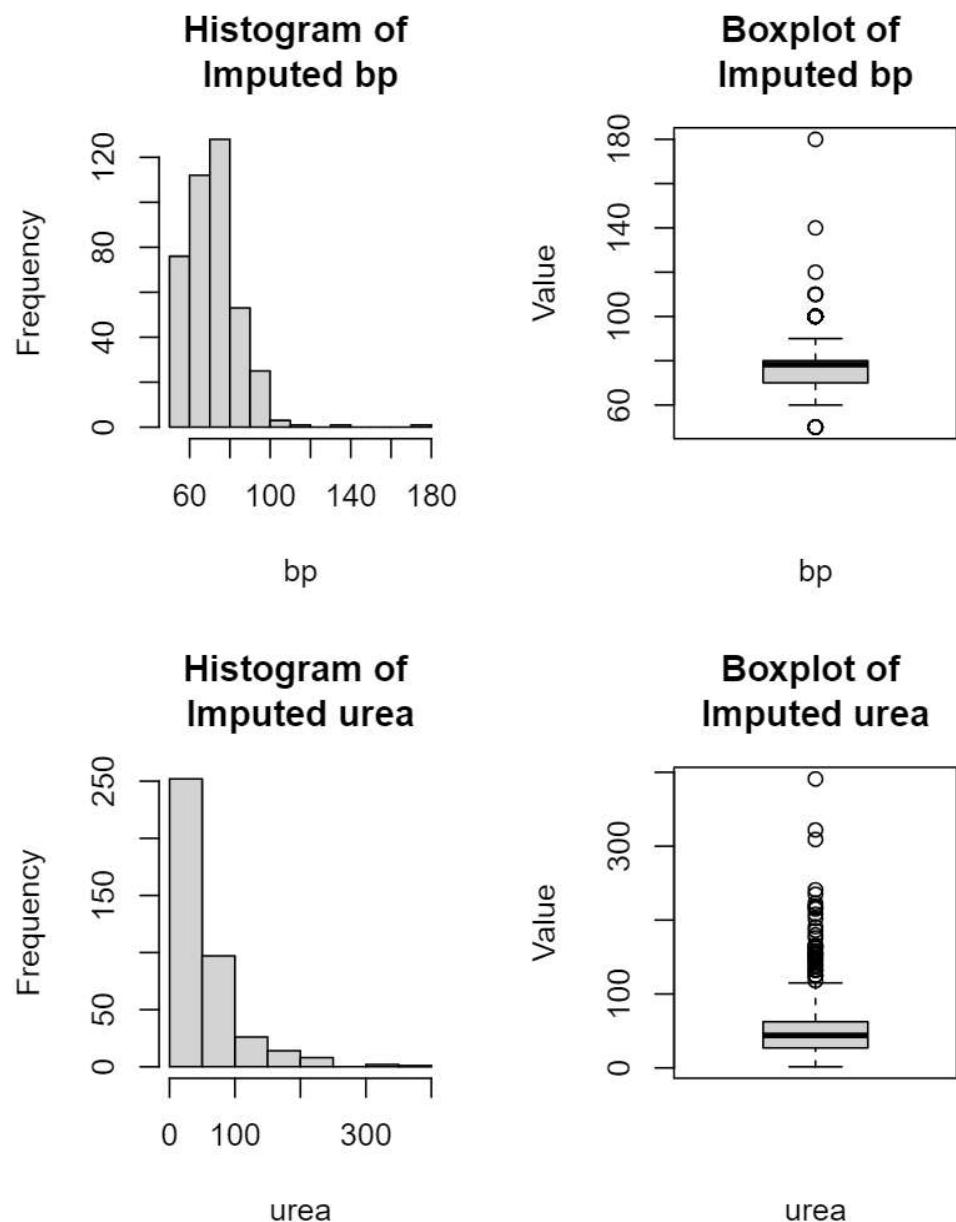
```

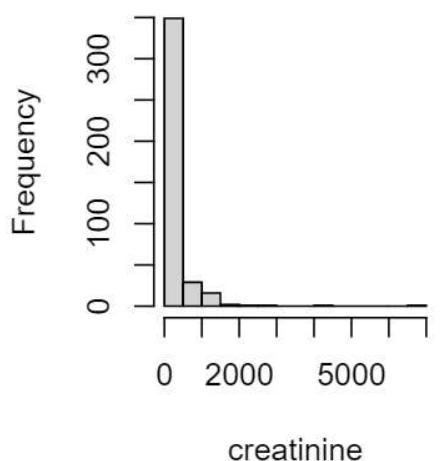
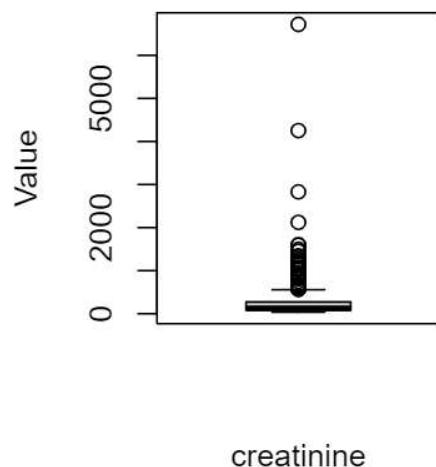
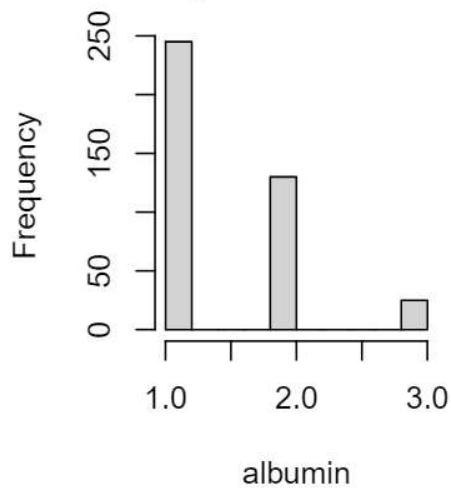
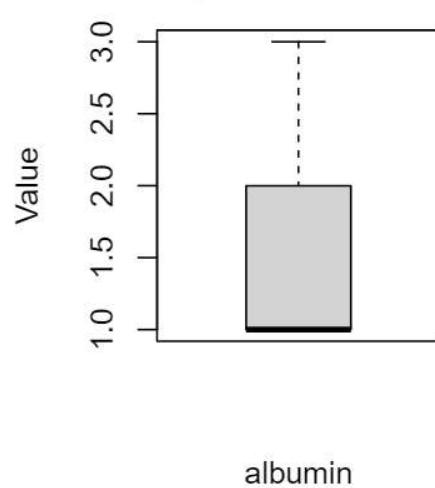
      age      bp      urea      creatinine      albumin      glucose
stats numeric,5  numeric,5  numeric,5  numeric,5  numeric,5  numeric,5
n      391       388       381       383       354       356
conf  numeric,2  numeric,2  numeric,2  numeric,2  numeric,2  numeric,2
out   numeric,10 numeric,36 numeric,38 numeric,51 numeric,0 numeric,34
group numeric,10 numeric,36 numeric,38 numeric,51 numeric,0 numeric,34
names ""
      "
```

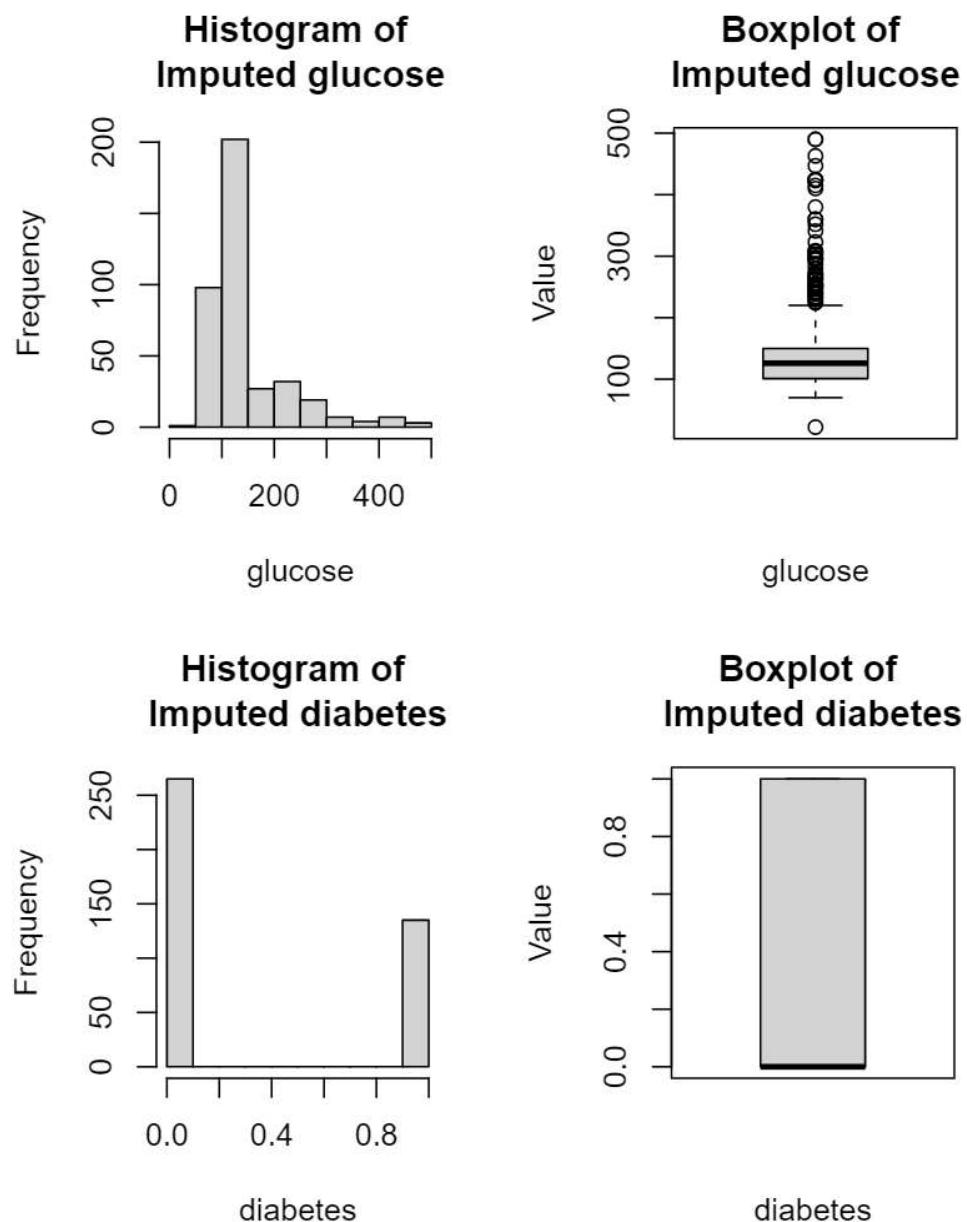
```
diabetes
stats numeric,5
n      400
conf   numeric,2
out    numeric,0
group  numeric,0
names  ""

1 # Store required imputed column names
2 numcols2 <- dx2 %>% colnames
3 # Use sapply() to apply sbsplots_i() to required columns.
4 sapply(numcols2, sbsplots_i, vars = data.frame(dx2)) # Imputed
```





Histogram of Imputed creatinine**Boxplot of Imputed creatinine****Histogram of Imputed albumin****Boxplot of Imputed albumin**



```

age          bp          urea        creatinine      albumin      glucose
stats numeric,5 numeric,5 numeric,5 numeric,5 numeric,5 numeric,5
n      400          400          400          400          400          400
conf numeric,2 numeric,2 numeric,2 numeric,2 numeric,2 numeric,2
out   numeric,10 numeric,36 numeric,39 numeric,44 numeric,0 numeric,53
group numeric,10 numeric,36 numeric,39 numeric,44 numeric,0 numeric,53
names   ""

```

```

        diabetes
stats numeric,5
n      400
conf   numeric,2
out    numeric,0
group  numeric,0
names  ""

1  ### Explanation:
2  # On looking at the boxplot for given imputed and non-imputed variables,
3  # the imputed variables appear to have a more symmetric distribution of
4  # values. This can be observed with Age for instance. After imputation
5  # with mean, the median of the Age values is much closer to the center
6  # whereas in case of the Non-Imputed Age, the median is much closer
7  # to the second quartile making it skewed to the right. This can be
8  # observed in case of blood pressure as well.
9
10 # In case of urea, the Non-imputed values were more skewed towards the
11 # left and on imputation, the Median has moved much closer to the centre.
12
13 # We might not see much of a difference incase of creatinine but we see a
14 # small reduction in outliers on imputation. Finally for the discrete
15 # variable albumin, there is no difference in the boxplot although
16 # the histogram shows an increase in "normo" values on imputation to mode.
17 # The variable 'Diabetes' does not see any change overall as there are
18 # a negligible amount of null values pre-imputation.

```

Problem 3.c (6 points)

- Plot a single figure containing boxplots of potential predictors for diabetes grouped by cases and controls. (Hint : `par(mfrow=c(1,5))`)
- Use these to decide which predictors to keep for future analysis.
- For any categorical variables create a table instead. Justify your answers.

```

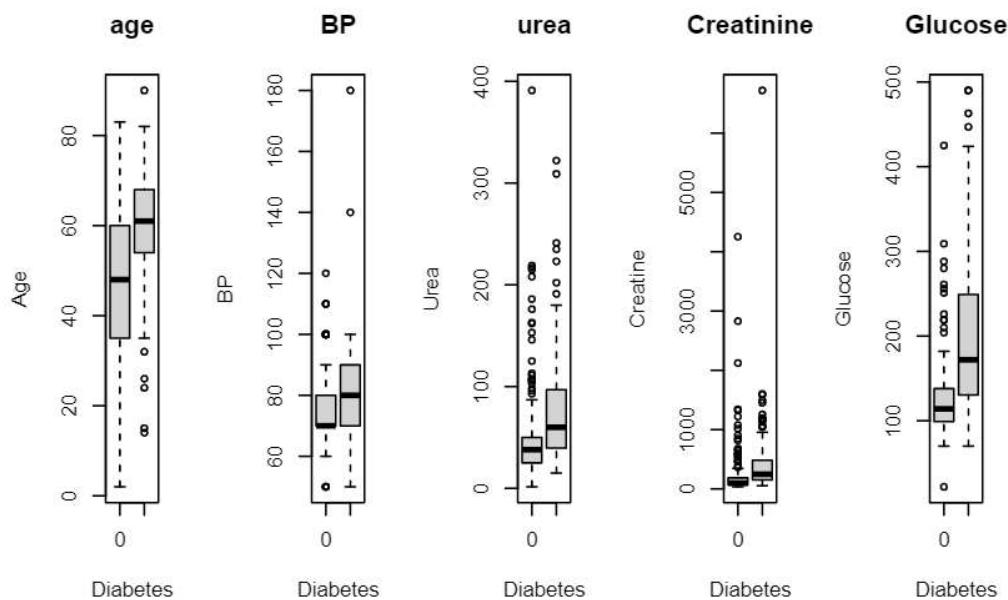
1  # Split page for 5 plots:
2  par(mfrow=c(1,5))
3  # Now to plot box plots of potential predictor variables
4  # (age,bp,urea,creatinine,glucose)
5  boxplot(dx2$age~dx2$diabetes,
6           main='age',xlab = 'Diabetes',ylab = 'Age')

```

```

7   boxplot(dx2$bp~dx2$diabetes,
8       main='BP',xlab = 'Diabetes',ylab = 'BP')
9   boxplot(dx2$urea~dx2$diabetes,
10      main='urea',xlab = 'Diabetes',ylab = 'Urea')
11  boxplot(dx2$creatinine~dx2$diabetes,
12      main='Creatinine',xlab = 'Diabetes',ylab = 'Creatine')
13  boxplot(dx2$glucose~dx2$diabetes,
14      main='Glucose',xlab = 'Diabetes',ylab = 'Glucose')

```



```

1 # Frequency table: albumin v diabetes.
2 # Rows indicate Albumin values and columns indicate diabetes values.
3 table(dx2$albumin,dx2$diabetes)

```

0	1
1	192 53
2	61 69
3	12 13

```

1 #### Justification:
2 # On observing the boxplots of several predictor variables in relation
3 # to diabetes, age and glucose appear to be more viable. Not only do they

```

```

4 # have the least number of outliers but the data points are also
5 # more symmetrically disturbed compared to other variables.
6 # For instance, there are too many outliers in creatinine and bp
7 # (diabetes =0) samples are skewed to the right.

```

Problem 3.d (9 points)

- Use your findings from the previous exercise and fit an appropriate model of diabetes with two predictors.
- Print a summary and explain the results as you would communicate it to a colleague with a medical background with a very little statistical knowledge.

```

1 # Based on previous findings, we plot a logistic regression model with
2 # chosen predictor variables using glm().
3 diabetes_model1<-glm(diabetes~age+glucose,data = dx2,family = 'binomial')
4 # Let's view a summary of the model.
5 summary(diabetes_model1)

```

Call:

```
glm(formula = diabetes ~ age + glucose, family = "binomial",
     data = dx2)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.7764	-0.7237	-0.4284	0.6580	2.7543

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-6.066415	0.661022	-9.177	< 2e-16 ***
age	0.049688	0.009481	5.241	1.60e-07 ***
glucose	0.018031	0.002532	7.120	1.08e-12 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

```
Null deviance: 511.49 on 399 degrees of freedom
Residual deviance: 368.02 on 397 degrees of freedom
AIC: 374.02
```

Number of Fisher Scoring iterations: 5

```

1  #### Summary Explanation:
2  # The first line indicates the logistic regression model with the parameters
3  # that were inputted.
4
5  # The second line contains "deviance residuals" that are simply a 5-number
6  # summary of the residuals (just like in box-plot). We have minimum and
7  # maximum values that are roughly the same distance from origin if you think
8  # about it. We also see quantiles 1 followed by median and quantile 3 at last.
9
10 # Third line has the coefficients! The intercepts are values when the
11 # features are at 0. The estimates give expected change in response due
12 # to a unit change in the feature. Next we have standard errors that are
13 # simply errors involved in calculating estimates that allow to construct
14 # confidence intervals for that feature.
15
16 # Fourth line has the Z test values that describe how far the estimated value
17 # value is from null hypothesis value.
18
19 # Next we have the "p-values" ( $p(>|z|)$ ). This is the probability of observing
20 # a value as extreme as z. We reject the null hypothesis if this value is
21 # very low. These probabilities follow Significance Codes below. The
22 # interpretation is:
23 # if  $p(>|z|)$  in (0,0.001]
24 # (0,0.001] then ***
25 # (0.001,0.1] then **
26 # (0.1,1] then *
27 # 0.05 is alpha value.
28 # Since all values are in (0,0.001], we see '***' next to every  $Pr(|>z|)$ value.
29
30 # Dispersion parameter indicates whether the distribution is wide or narrow.
31 # For the binomial family(as specified in input), this value is 1.
32
33 # Next we have the values, 'Null deviance' and "Residual deviance". There is a
34 # significant difference between the values(143.47) which is an indication that
35 # the model is a good fit.
36
37 # The "Akaike information criterion (AIC)" is another measure of model fit.The
38 # lower its value compared to other models, the better the fit. Although

```

```
39 # individually, it is not that significant.  
40  
41 # Finally, fisher scoring iterations (5) is just a measure of how long it took  
42 # to fit the model.
```

Problem 4 (19 points)

Problem 4.a. (9 points)

- Add a third predictor to the final model from problem 3, perform a likelihood ratio test to compare both models and report the p-value for the test.
- Is there any support for the additional term?
- Plot a ROC curve for both models and report the AUC, explain the results as you would communicate it to a colleague with a medical background with a very little statistical knowledge.
- Print a summary and explain the results as you would communicate it to a colleague with a medical background with a very little statistical knowledge.

```

1 # Add a 3rd predictor (urea) to previous model.
2 diabetes_model2<-glm(diabetes~age+glucose+urea,
3                         data = dx2,family = 'binomial')
4 # Model 1:
5 # Odds ratio for age and confidence interval
6 or.age <- exp(coef(diabetes_model1)[2])
7 ci.age <- exp(confint(diabetes_model1))[2, ]

```

Waiting for profiling to be done...

```

1 # Round Odds ratio for age and confidence interval to 2 decimal places
2 round(c(or.age, ci.age), 2)

```

	age	2.5 %	97.5 %
1.05	1.03	1.07	

```

1 # Odds ratio for glucose and confidence interval
2 or.glucose <- exp(coef(diabetes_model1)[3])
3 ci.glucose <- exp(confint(diabetes_model1))[3, ]

```

Waiting for profiling to be done...

```
1 # Round Odds ratio and confidence interval to 2 decimal places
2 round(c(or.glucose, ci.glucose), 2)
```

glucose 2.5 % 97.5 %
1.02 1.01 1.02

```
1 # Model 2:
2 # Odds ratio for age and confidence interval
3 or.age <- exp(coef(diabetes_model2)[2])
4 ci.age <- exp(confint(diabetes_model2))[2, ]
```

Waiting for profiling to be done...

```
1 # Round Odds ratio and confidence interval to 2 decimal places
2 round(c(or.age, ci.age), 2)
```

age 2.5 % 97.5 %
1.05 1.03 1.07

```
1 # Odds ratio for glucose and confidence interval
2 or.glucose <- exp(coef(diabetes_model2)[3])
3 ci.glucose <- exp(confint(diabetes_model2))[3, ]
```

Waiting for profiling to be done...

```
1 # Round Odds ratio and confidence interval to 2 decimal places
2 round(c(or.glucose, ci.glucose), 2)
```

glucose 2.5 % 97.5 %
1.02 1.01 1.02

```
1 # Odds ratio for urea and confidence interval
2 or.urea <- exp(coef(diabetes_model2)[4])
```

```
3 ci.urea <- exp(confint(diabetes_model2))[4, ]
```

Waiting for profiling to be done...

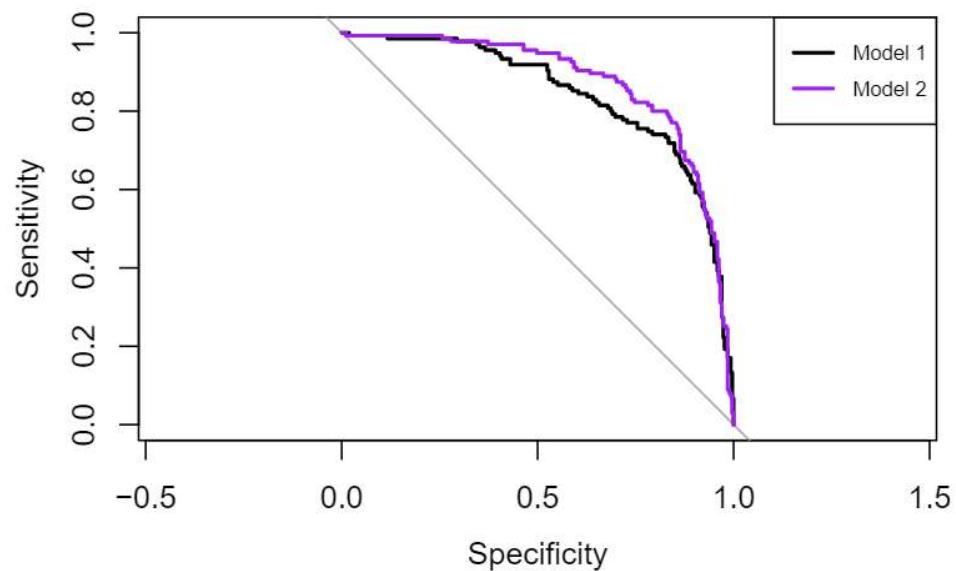
```
1 # Round Odds ratio and confidence interval to 2 decimal places
2 round(c(or.urea, ci.urea), 2)
```

	urea	2.5 %	97.5 %
1	1.01	1.01	1.02

```
1 # Report the p value using pchisq()
2 p<-pchisq(diabetes_model1$deviance-diabetes_model2$deviance,
3             df=1,lower.tail = FALSE)
4 cat('The p value is',p)
```

The p value is 3.822726e-06

```
1 # Plot a roc curve to compare both models.
2 suppressMessages(invisible({
3   roc(dx2$diabetes, diabetes_model1$fitted.values,
4       plot = TRUE, xlim = c(0,1))
5   # add = TRUE adds the roc curve to the previous plot
6   roc(dx2$diabetes, diabetes_model2$fitted.values,
7       plot = TRUE, xlim = c(0,1),
8       add = TRUE, col = "purple")
9   legend('topright',legend=c('Model 1','Model 2'),
10         col=c('black','purple'),
11         lty=c(1,1),lwd=c(2,2),cex = 0.7)
12 }))
```



```

1 # Print area under roc curve
2 cat('Area under roc curve',
3     roc(dx2$diabetes, diabetes_model2$fitted.values,
4     )$auc, '\n')

```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

Area under roc curve 0.8745073

```

1 #### ROC Curve explanation:
2 # The area under the roc curve is a good measure of the performance
3 # of a model. It is simply a plot between the False Positive Rate
4 # (Specificity) and True Positive Rate (Sensitivity). It is used
5 # to test performance of classification models in Machine Learning.
6 # It can also be used to test how significant the addition of a
7 # predictor variable is to the performance of model.
8 # As you can see below, there are 2 roc curves. The black one is
9 # that of model 1 and the purple one is that of model 2. Model 2
10 # is model 1 with an additional predictor variable (urea). Since
11 # auc did improve, urea is indeed a valid addition and there is

```

```

12 # clearly support to it.
13
14 # Examine summary of the optimum model.
15 summary(diabetes_model2)

```

Call:

```
glm(formula = diabetes ~ age + glucose + urea, family = "binomial",
  data = dx2)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-3.0576	-0.6495	-0.3912	0.6083	2.9107

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-6.590163	0.707935	-9.309	< 2e-16 ***
age	0.047483	0.009922	4.786	1.70e-06 ***
glucose	0.017042	0.002472	6.894	5.42e-12 ***
urea	0.012717	0.002959	4.298	1.73e-05 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 511.49 on 399 degrees of freedom
 Residual deviance: 346.67 on 396 degrees of freedom
 AIC: 354.67

Number of Fisher Scoring iterations: 5

```

1 ### Summary Explanation:
2 # The first line indicates the logistic regression model with the parameters
3 # that were inputted.
4
5 # The second line contains "deviance residuals" that are simply a 5-number
6 # summary of the residuals (just like in box-plot). We have minimum and
7 # maximum values that are roughly the same distance from origin if you think
8 # about it. We also see quantiles 1 followed by median and quantile 3 at last.
9

```

```

10 # Third line has the coefficients! The intercepts are values when the
11 # features are at 0. The estimates give expected change in response due
12 # to a unit change in the feature. Next we have standard errors that are
13 # simply errors involved in calculating estimates that allow to construct
14 # confidence intervals for that feature.
15
16 # Fourth line has the Z test values that describe how far the estimated value
17 # value is from null hypothesis value.
18
19 # Next we have the "p-values" ( $p(>|z|)$ ). This is the probability of observing
20 # a value as extreme as z. We reject the null hypothesis if this value is
21 # very low. These probabilities follow Significance Codes below. The
22 # interpretation is:
23 # if  $p(>|z|)$  in (0,0.001]
24 # (0,0.001] then ***
25 # (0.001,0.1] then **
26 # (0.1,1] then *
27 # 0.05 is alpha value.
28 # Since all values are in (0,0.001], we see '***' next to every  $Pr(|>z|)$  value.
29
30 # Dispersion parameter indicates whether the distribution is wide or narrow.
31 # For the binomial family(as specified in input), this value is 1.
32
33 # Next we have the values, 'Null deviance' and "Residual deviance". There is a
34 # significant difference between the values(164.82) which is an indication that
35 # the model is a good fit.
36
37 # The "Akaike information criterion (AIC)" is another measure of model fit.The
38 # lower its value compared to other models, the better the fit. Although
39 # individually, it is not that significant.
40
41 # Finally, fisher scoring iterations (5) is just a measure of how long it took
42 # to fit the model.

```

Problem 4.b (10 points)

- Perform 10-folds cross validation for your chosen model based on the above answers.
- Report the mean cross-validated AUCs in 3 significant figures.

```

1  # Initialize number of cross validation folds.
2  num.folds<-10
3  # Next we create 10 cross validation folds using the
4  # createFolds() function.
5  folds <- createFolds(dx2$diabetes, k = num.folds)
6  # Initialize predictor variable.
7  pred.cv <- NULL
8  # Train object to store observed training values.
9  regr.cv<-NULL
10 auc.cv <- numeric(num.folds)
11 for(f in 1:num.folds){
12   # For each fold:
13   # Store the test values
14   test.idx <- folds[[f]]
15   # Use setdiff() to exclude the k-fold rows from
16   # original data and thus form train data
17   train.idx <- setdiff(1:nrow(dx2), folds[[f]])
18   # The regr.cv is the logistic regression object
19   regr.cv[[f]] <- glm(diabetes~age+glucose+urea,
20                      data = dx2,subset = train.idx,
21                      family = 'binomial')
22   # Now we use the predict() function to get prediction
23   # values.
24   pred.cv[[f]] <- data.frame(obs = dx2$diabetes[test.idx]
25                               ,pred = predict(regr.cv[[f]] ,
26                               newdata = dx2,
27                               type = "response"))[test.idx])
28   # Store area under curve for the fold
29   auc.cv[f] <- roc(obs ~ pred, data = pred.cv[[f]])$auc
30
31 }

```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

Setting levels: control = 0, case = 1

Setting direction: controls < cases

Setting levels: control = 0, case = 1

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

```
1 # Mean cross-validated AUC is simply mean of auc.cv vector.  
2 # Round resultan to 3 decimal places:  
3 cat('The mean of aucs is',round(mean(auc.cv),3))
```

```
The mean of aucs is 0.867
```