

Assignment 1

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

Aysu Ismayilova

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 library(data.table)
2 library(dplyr) #dplyr loads magrittr too
3 longegfr1.dt <- fread("data_assignment1/longegfr1.csv", stringsAsFactors = F, key = 'id')
4
5 longegfr2.dt <- fread("data_assignment1/longegfr2.csv", stringsAsFactors = F)
6
7 colnames(longegfr2.dt) <- c("id", "fu.years", "egfr") #renaming column names so
8 #they are consistent in the two datasets
9
10 #combining two dataframes
11 longegfr.dt <- merge(longegfr1.dt, longegfr2.dt, by = c("id", "fu.years"), all = TRUE)
12
13
14
15
16 # Order the merged data table by subject identifier and follow-up time
17 longegfr_ordered_head_10.dt <- head(longegfr.dt %>% arrange(id, fu.years), n=10)
18 kable(longegfr_ordered_head_10.dt, caption='longegfr data')|>kable_styling(full_width=F,
19 position='center',latex_options = 'basic')
```

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()`.

Table 1: longegfr data

id	fu.years	sex	baseline.age	egfr
1	0.0000	0	65.5	76.48
1	0.1533	0	65.5	47.36
1	0.6899	0	65.5	94.87
1	1.1882	0	65.5	52.12
1	1.8398	0	65.5	91.91
1	2.2806	0	65.5	76.52
1	3.3895	0	65.5	46.79
1	3.7563	0	65.5	35.56
1	4.5229	0	65.5	28.41
1	5.3607	0	65.5	20.85

- Tabulate the number of patients with average eGFR in the following ranges: (0, 15], (15, 30], (30, 60], (60, 90], (90, max(eGFR)).
- Count and report the number of patients with missing average eGFR.

```

1  #Answer in this chunk
2  library(dplyr)
3
4
5
6
7
8  # compute the average eGFR and length of follow-up for each patient
9  longegfr_summary.dt <- longegfr.dt[, .(avg_egfr = mean(egfr, na.rm = TRUE),
10                                     follow_up_length = max(fu.years, na.rm = TRUE),
11                                     sex = first(sex),
12                                     age = first(baseline.age)), by = id]
13
14 # print the first 10 rows of the new dataset
15 longegfr_summary_10 <- head(longegfr_summary.dt, n = 10)
16 kable(longegfr_summary_10, caption='longegfr data')|>kable_styling(full_width=F,
17                                     position = 'center',
18                                     latex_options = 'basic')

1  longegfr_summary_no_na <- na.omit(longegfr_summary.dt)
2  # tabulate the number of patients with average eGFR in each range
3  bin_edges <- c(0, 15, 30, 60, 90, max(longegfr_summary_no_na$avg_egfr, na.rm = TRUE))

```

Table 2: longegfr data

id	avg_eGFR	follow_up_length	sex	age
1	43.04333	6.4586	0	65.5
2	38.93294	2.0698	1	83.2
3	85.72000	6.5161	1	19.6
4	76.59308	5.2786	0	50.3
5	13.90892	6.3929	1	72.1
6	85.66435	6.2313	1	65.5
7	64.21758	5.8453	0	73.0
8	66.28333	1.5606	1	84.7
9	86.35750	5.8700	0	72.6
10	107.00429	5.1964	0	50.4

Table 3: longegfr data

Breaks	Number of patients
(0, 15]	2
(15, 30]	9
(30, 60]	84
(60, 90]	86
(90, max(eGFR)]	66

```

4 num_patients <- cut(longegfr_summary_no_na$avg_eGFR,
5 breaks = c(-Inf, 15, 30, 60, 90, max(longegfr_summary_no_na$avg_eGFR)),
6 labels = c("(0, 15]", "(15, 30]", "(30, 60]", "(60, 90]",
7 "(90, max(eGFR)]"))
8 n_pat <- table(num_patients)
9 eGFR_bins <- cut(longegfr_summary.dt$avg_eGFR, breaks = bin_edges,
10 include.lowest = TRUE)
11 kable(n_pat, caption='longegfr data', col.names = c('Breaks',
12 'Number of patients')) |> kable_styling(full_width=F,
13 position = 'center',
14 latex_options = 'basic')

```

```

1 # count and report the number of patients with missing average eGFR
2 sum(is.na(longegfr_summary.dt$avg_eGFR))

```

[1] 3

Table 4: Data with patients with average eGFR in the (90,max(eGFR))

id	sex	age	avg_eGFR	follow_up_length
10	0	50.4	107.00429	5.1964
14	0	65.1	116.09200	4.0986
25	0	40.1	95.35625	4.2847
31	0	74.8	113.59250	1.4675
32	0	67.6	126.98950	6.2177
33	0	74.2	116.35000	1.6016
45	1	24.9	91.25000	0.0000
49	1	68.2	128.25800	6.1602
52	1	56.3	93.31544	6.4805
56	1	56.2	129.70909	5.9630

There are 3 patients with missing eGFR

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 library(data.table)
2 num_high_eGFR <- nrow(subset(longegfr_summary.dt, avg_eGFR > 90))
3 #
4 # # Filter the summary table to include only patients with avg_eGFR > 90
5 longegfr_high.dt <- subset(longegfr_summary.dt, avg_eGFR > 90,
6 select = c(id, sex, age, avg_eGFR, follow_up_length))
7 longegfr_high.dt_10 <- head(longegfr_high.dt, n=10)
8 kable(longegfr_high.dt_10, caption='Data with patients
9 with average eGFR
10 in the (90,max(eGFR))')|>kable_styling(full_width=F,position ='center',
11 latex_options = 'basic')

1 # Print summary of the new data table
2 summary(longegfr_high.dt)

```

	id	sex	age	avg_eGFR
Min.	: 10.00	Min. :0.0000	Min. :22.10	Min. : 90.04

1st Qu.: 86.25	1st Qu.:0.0000	1st Qu.:47.20	1st Qu.: 99.13
Median :144.00	Median :0.0000	Median :55.20	Median :109.81
Mean :141.88	Mean :0.3333	Mean :55.27	Mean :112.13
3rd Qu.:197.50	3rd Qu.:1.0000	3rd Qu.:63.80	3rd Qu.:123.20
Max. :250.00	Max. :1.0000	Max. :90.90	Max. :147.69

follow_up_length

Min. :0.000

1st Qu.:1.671

Median :4.857

Mean :4.064

3rd Qu.:5.937

Max. :6.590

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 library(ggplot2)
2 library(dplyr)
3
4 # Filter data for patients 3, 37, 162, and 223
5 patient_ids <- c(3, 37, 162, 223)
6 patient_data <- longegfr.dt %>% filter(id %in% patient_ids)
7
8 # Create a function to compute the linear regression and confidence intervals
9 compute_regression <- function(data) {
10   fit <- lm(egfr ~ fu.years, data = data)
11   conf_int <- confint(fit, level = 0.95)
12   data.frame(intercept = fit$coefficients[1],
13              slope = fit$coefficients[2],
14              lower = conf_int[1, ],

```

```

15         upper = conf_int[2, ]))
16     }
17
18     # Compute the regression lines and confidence intervals for each patient
19     regression_data <- patient_data %>% group_by(id) %>% do(compute_regression(.))

```

Warning in data.frame(intercept = fit\$coefficients[1], slope = fit\$coefficients[2], : row names were found from a short variable and have been discarded

Warning in data.frame(intercept = fit\$coefficients[1], slope = fit\$coefficients[2], : row names were found from a short variable and have been discarded

Warning in data.frame(intercept = fit\$coefficients[1], slope = fit\$coefficients[2], : row names were found from a short variable and have been discarded

Warning in data.frame(intercept = fit\$coefficients[1], slope = fit\$coefficients[2], : row names were found from a short variable and have been discarded

```

1  # Remove extreme eGFR values for each patient and compute regression lines again
2  filtered_data <- patient_data %>% group_by(id) %>% filter(egfr != min(egfr) &
3  egfr != max(egfr))
4  filtered_regression_data <- filtered_data %>% group_by(id) %>% do(compute_regression(.))

```

Warning in data.frame(intercept = fit\$coefficients[1], slope = fit\$coefficients[2], : row names were found from a short variable and have been discarded

Warning in data.frame(intercept = fit\$coefficients[1], slope = fit\$coefficients[2], : row names were found from a short variable and have been discarded

Warning in data.frame(intercept = fit\$coefficients[1], slope = fit\$coefficients[2], : row names were found from a short variable and have been discarded

```

1 # Combine the two sets of
2 #regression lines and confidence
3 #intervals into a single data frame
4 all_regression_data <- rbind(regression_data, filtered_regression_data)
5
6 # Create the plot
7 ggplot(patient_data, aes(x = fu.years, y = egfr, color = as.factor(id))) +
8   geom_point() +
9   geom_smooth(method = "lm", se = FALSE,
10 aes(group = id)) +
11   geom_smooth(data = filtered_data, method = "lm", se = FALSE,
12 aes(group = id, color = "Filtered"), linetype = "dashed") +
13   geom_hline(yintercept = 60, linetype = "dotted") +
14   facet_wrap(~id) +
15   scale_color_discrete(name = "Patient") +
16   labs(x = "Follow-up Time (Years)", y = "eGFR") +
17   theme_bw()

```

`geom_smooth()` using formula = 'y ~ x'

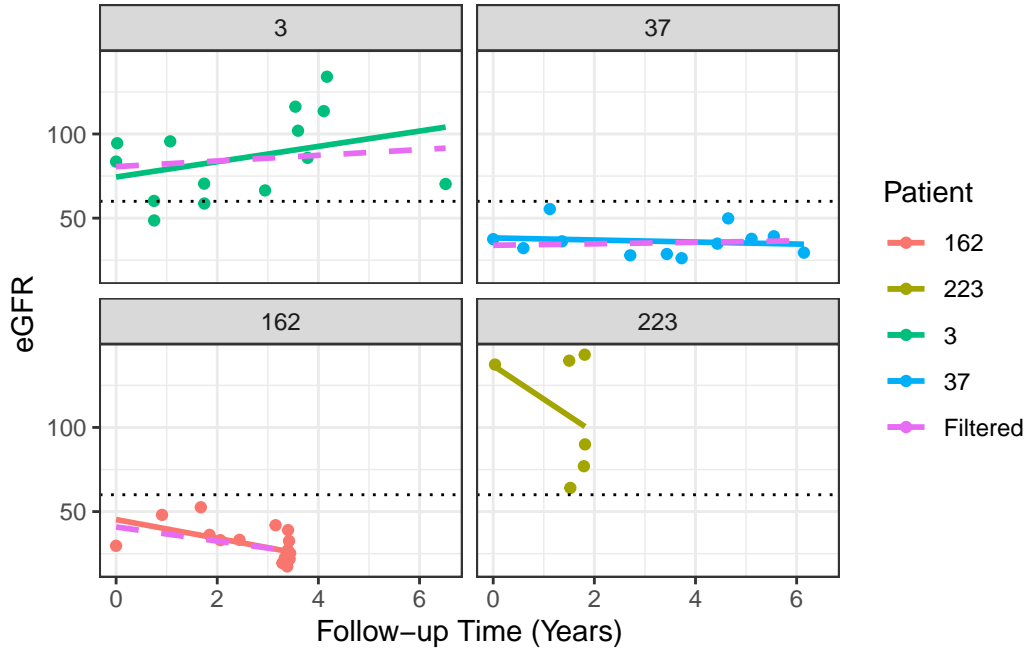
Warning: Removed 1 rows containing non-finite values (`stat_smooth()`).

`geom_smooth()` using formula = 'y ~ x'

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

Table 5: 95regression coefficients of the fitted model.

	2.5 %	97.5 %
(Intercept)	49.57033	87.3760083
fu.years	-11.37162	0.8066081



```

1 model <- lm(egfr ~ fu.years, data = patient_data)
2 confidence_intervals <- confint(model, level = 0.95)
3 kable(confidence_intervals, caption='95% confidence interval for the
4   regression coefficients of the fitted model.')|>kable_styling(full_width=F,
5
6   position = 'center', latex_options = 'basic')

```

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 library(data.table)
2 library(dplyr) #dplyr loads magrittr too
3
4 scr.dt <- fread("data_assignment1/scr.csv", stringsAsFactors = FALSE)
5 summary_of_scr <- summary(scr.dt$scr)
6
7 iqr <- IQR(scr.dt$scr, na.rm = TRUE)
8 vec <- c(summary_of_scr, iqr)
9 kable(vec,
10 caption='inter-quartile
11 range of serum
12 creatinine')|>kable_styling(full_width=F,
13 position = 'center',
14 latex_options = 'basic')
```

Table 6: inter-quartile range of serum creatinine

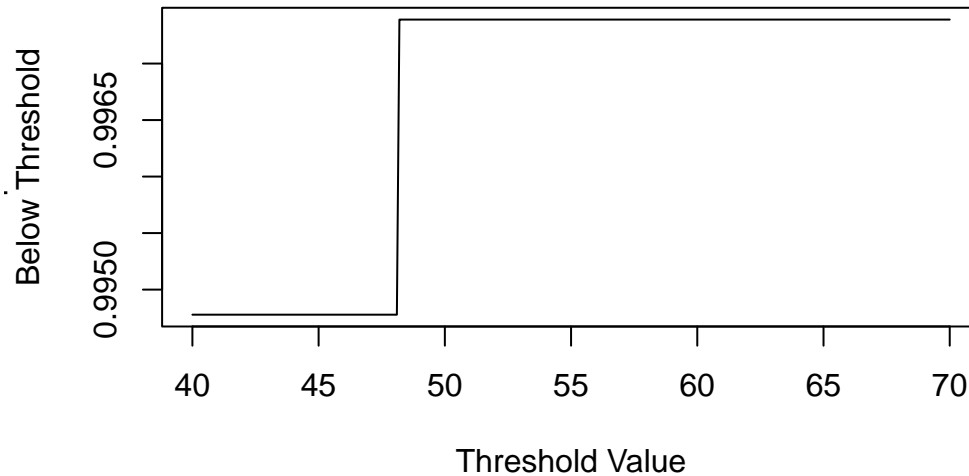
	x
Min.	0.400000
1st Qu.	0.900000
Median	1.300000
Mean	3.072454
3rd Qu.	2.800000
Max.	76.000000
NA's	18.000000
	1.900000

- The MDRD4 equation estimates the glomerular filtration rate (GFR) based on serum creatinine, age, gender, and race. It was developed using data from a study of patients with chronic kidney disease (CKD). The CKD-EPI equation is an improvement over the MDRD4 equation, also estimating GFR based on serum creatinine, age, gender, and race, but with different coefficients. The CKD-EPI equation is considered more accurate than the MDRD4 equation, particularly for higher GFR levels. It was developed using data from multiple studies involving patients with CKD.

```

1  # Generate a range of values to test as the threshold
2  thresholds <- seq(40, 70, by = 0.1)
3
4  # Compute the proportion of values below each threshold
5  proportions <- sapply(thresholds, function(t) mean(scr.dt$scr < t, na.rm=TRUE))
6
7  # Plot the proportion of values below each threshold
8  plot(thresholds, proportions, type = "l", xlab = "Threshold Value",
9       ylab = "Proportion
10      Below Threshold")

```



```

1 # Identify the threshold value that maximizes
2 #the proportion of values below the threshold
3 max_threshold <- thresholds[which.max(proportions)]

```

The typical range for serum creatinine is: **For adult men, 0.74 to 1.35 mg/dL (65.4 to 119.3 micromoles/L) For adult women, 0.59 to 1.04 mg/dL (52.2 to 91.9 micromoles/L)**

A range of threshold values from 40 to 70 was chosen with an increment of 0.1 for finding the optimal threshold and confirm this. The mean function was used to calculate the proportion of values that are below each threshold, and sapply function was used to apply this function to every threshold value in the vector. Then, a plot was created between the proportion of values below each threshold and the threshold values using the plot function. The threshold value that gives the highest proportion of values below the threshold was determined using the which.max function, and its value was saved in the max_threshold variable.

The conversion factor from $\mu\text{mol/L}$ to mg/dL for serum creatinine is 0.0113, which can be obtained using the formula provided in the text or by consulting standard conversion tables. In order to convert the serum creatinine values from $\mu\text{mol/L}$ to mg/dL , we need to divide them by 88.42.

After making any changes to the serum creatinine values, it is important to assess the distribution of the modified data. This can be done using summary statistics such as mean, median, standard deviation, and range to compare the original and modified data. Additionally, we created histograms and boxplots to visualize the distributions and identify any outliers or other patterns. It is important to ensure that the changes made did not significantly alter the data.

```

1 # Calculate the proportion of missing values
2 prop_missing <- sum(is.na(scr.dt$scr)) / nrow(scr.dt)
3 cat(paste0("Proportion of missing values: ", round(prop_missing, 2), "\n"))

```

Proportion of missing values: 0.04

```

1 # Convert  $\mu\text{mol/L}$  to mg/dL
2 scr.dt_new <- ifelse(scr.dt$scr > 52.2, scr.dt$scr / 88.42, scr.dt$scr)
3
4 # Summarize the data
5 summary(scr.dt_new)

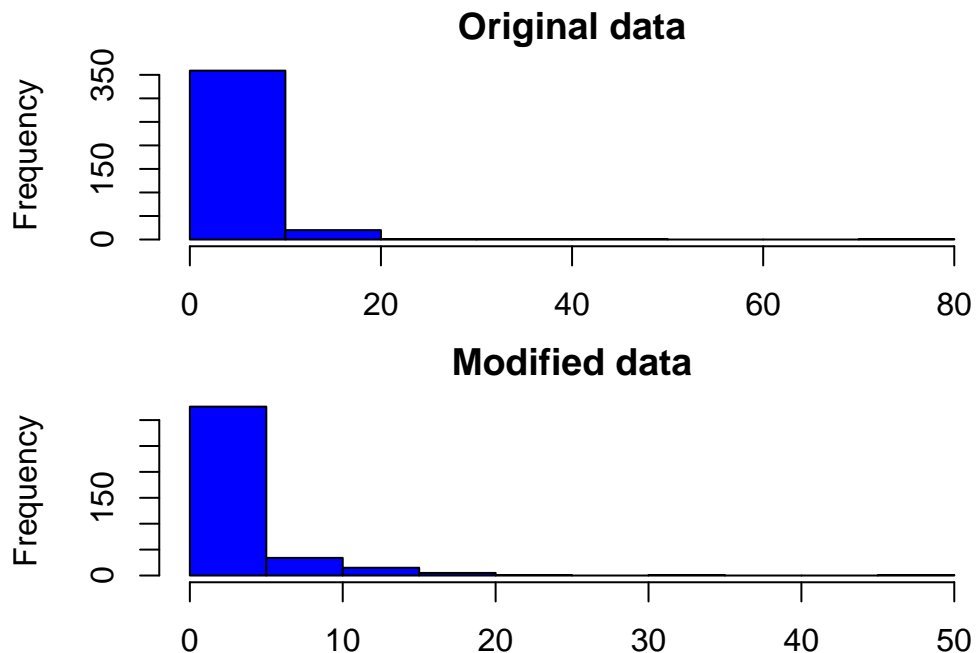
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0.400	0.900	1.300	2.876	2.800	48.100	18

```

1 # Create histograms of the original and modified data
2 par(mfrow = c(2, 1))
3 par(mar = c(2, 4, 2, 2))
4 hist(scr.dt$scr, main = "Original data", xlab = "Serum creatinine (mg/dL)", col = "blue")
5 hist(scr.dt_new, main = "Modified data", xlab = "Serum creatinine (mg/dL)", col = "blue")

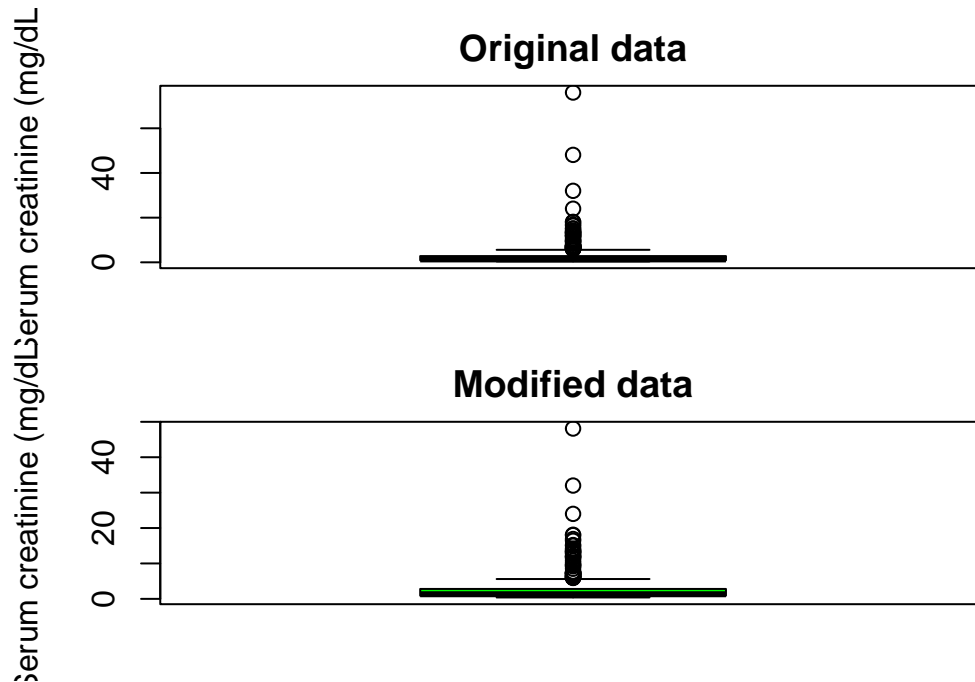
```



```

1 # Create box plots of the original and modified data
2 par(mfrow = c(2, 1))
3 par(mar = c(2, 4, 2, 2))
4 boxplot(scr.dt$scr, main = "Original data", ylab = "Serum creatinine (mg/dL)",
5 col = "green")
6
7 boxplot(scr.dt_new, main = "Modified data", ylab = "Serum creatinine (mg/dL)",
8 col = "green")

```



```

1 summary_of_scr <- summary(scr.dt_new)
2 iqr <- IQR(scr.dt_new, na.rm = TRUE)
3 vec <- c(summary_of_scr, iqr)
4 kable(vec, caption='inter-quartile range of serum creatinine')|>kable_styling(full_width=F,
5 position = 'center', latex_options = 'ba

```

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.

Table 7: inter-quartile range of serum creatinine

	x
Min.	0.400000
1st Qu.	0.900000
Median	1.300000
Mean	2.876265
3rd Qu.	2.800000
Max.	48.100000
NA's	18.000000
	1.900000

- 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  scr_no_na <- na.omit(scr.dt)
2
3  # Calculate eGFR using MDRD4 equation
4  mdrd4 <- with(scr_no_na, 175 * (scr)^(-1.154) * age^(-0.203) * ifelse(sex == "Female"
5  0.742, 1) * ifelse(ethnic == "Black", 1.212, 1))
6
7  # Calculate eGFR using CKD-EPI equation
8  kappa <- ifelse(scr_no_na$sex == "Female", 0.7, 0.9)
9  alpha <- ifelse(scr_no_na$sex == "Female", -0.329, -0.411)
10 ckd_epi <- with(scr_no_na, 141 * (pmin(scr/ kappa, 1))^alpha *
11 (pmax(scr / kappa, 1))^(-1.209) * 0.993^age * ifelse(sex == "Female", 1.018, 1) *
12 ifelse(ethnic == "Black", 1.159, 1))
13
14 # Combine results into a data frame
15 egfr <- data.frame(mdrd4 = mdrd4, ckd_epi = ckd_epi)
16
17 # Print mean and standard deviation of MDRD4 and CKD-EPI
18 mdrd4_mean <- mean(mdrd4)
19 mdrd4_std <- sd(mdrd4)
20 ckd_epi_mean <- mean(ckd_epi)
21 ckd_epi_std <- sd(ckd_epi)
22
23 cat("Mean of MDRD4: ", round(mdrd4_mean, 2), "\n")

```

Mean of MDRD4: 59.34

```
1 cat("Standard Deviation of MDRD4: ", round(mdrd4_std, 2), "\n")
```

Standard Deviation of MDRD4: 47.67

```
1 cat("Mean of CKD-EPI: ", round(ckd_epi_mean, 2), "\n")
```

Mean of CKD-EPI: 58.5

```
1 cat("Standard Deviation of CKD-EPI: ", round(ckd_epi_std, 2), "\n")
```

Standard Deviation of CKD-EPI: 42.12

```
1 ckd_mdrd4_cor <- cor(egfr$ckd_epi, egfr$mdrd4)
```

```
2
```

```
3 cat("Pearson Correlation Coefficient between CKD-EPI and MDRD4:", ckd_mdrd4_cor)
```

Pearson Correlation Coefficient between CKD-EPI and MDRD4: 0.9714988

```
1 # Cut MDRD4 eGFR values into three strata
```

```
2
```

```
3 mdrd4_strata <- cut(egfr$mdrd4, breaks = c(0, 60, 90, Inf), labels = c("(0-60)",  
4 "(60-90)", "> 90))")
```

```
5
```

```
6
```

```
7 # Compute mean and standard deviation of CKD-EPI eGFR for each stratum
```

```
8 ckd_epi_mean_str <- tapply(egfr$ckd_epi, mdrd4_strata, mean, na.rm = TRUE)
```

```
9 ckd_epi_std_str <- tapply(egfr$ckd_epi, mdrd4_strata, sd)
```

```
10
```

```
11 # Print results
```

```
12 cat("Mean of CKD-EPI eGFR by MDRD4 stratum:\n")
```

Mean of CKD-EPI eGFR by MDRD4 stratum:

```
1 print(ckd_epi_mean_str)
```

(0-60)	(60-90)	(> 90)
26.55813	80.25483	119.87214

```
1 cat("Standard Deviation of CKD-EPI eGFR by MDRD4 stratum:\n")
```


Table 8: Mean and Standard deviation of CKD-EPI based on three different strata

	ckd_epi_mean_str	ckd_epi_std_str
(0-60)	26.55813	18.84544
(60-90)	80.25483	10.45372
(> 90)	119.87214	19.33130

Standard Deviation of CKD-EPI eGFR by MDRD4 stratum:

```

1 print(ckd_epi_std_str)

      (0-60) (60-90) (> 90)
18.84544 10.45372 19.33130

1 df <- data.frame(ckd_epi_mean_str,ckd_epi_std_str)
2 kable(df,caption='Mean and Standard deviation of CKD-EPI based on three different
3 strata')|>
4 kable_styling(full_width=F,position = 'center',latex_options = 'basic')

1 egfr_head_15 <- head(egfr,n=15)
2
3 kable(egfr_head_15,caption='15 values
4 of CKD-EPI
5 and MDRD4')|>kable_styling(full_width=F,
6                               position = 'center',latex_options = 'basic')

```

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 # create a scatter plot of the two eGFR vectors
2 plot(mdrd4, ckd_epi, xlab = "eGFR_mdrd", ylab = "eGFR_ckd_epi", main = "Relationship
3 between
4 eGFR equations")
5
6 # add vertical line for median
7 abline(v = median(mdrd4), col = "red", lty = 2)

```

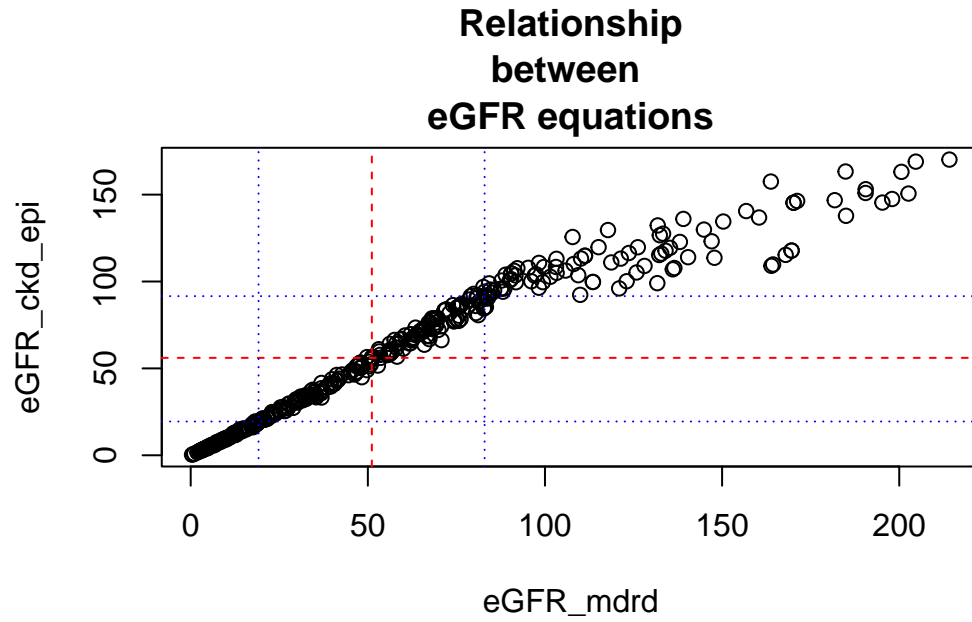
Table 9: 15 values of CKD-EPI and MDRD4

mdrd4	ckd_epi
47.948478	53.397906
184.850199	163.294281
12.678846	13.252444
53.428078	57.761862
68.281993	72.578746
1.897907	1.651094
99.676006	108.322818
27.759499	29.787787
8.010292	7.864905
11.851513	12.281808
23.987121	23.998393
23.420792	23.587189
12.770744	12.166705
17.676195	16.205967
6.085256	6.085585

```

8  abline(h = median(ckd_epi), col = "red", lty = 2)
9
10 # add vertical lines for first and third quantiles
11 abline(v = quantile(mdrd4, 0.25), col = "blue", lty = 3)
12 abline(v = quantile(mdrd4, 0.75), col = "blue", lty = 3)
13 abline(h = quantile(ckd_epi, 0.25), col = "blue", lty = 3)
14 abline(h = quantile(ckd_epi, 0.75), col = "blue", lty = 3)

```



Based on the scatter plot, the relationship between the two eGFR equations appears to be approximately linear. The points on the scatter plot are distributed somewhat evenly around the diagonal line, which suggests that the differences between the two eGFR values are proportional and not dependent on the value itself. Additionally, the vertical and horizontal lines corresponding to the median, first, and third quantiles are almost parallel, which further supports the idea of a linear 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 on 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()`.

Table 10: Cohort data

id	yob	age	bp	diabetes	albumin	urea	creatinine	glucose
PID_1	1971	48	80	1	micro	36	106.104	121
PID_2	2012	7	50	0	macro	18	70.736	NA
PID_3	1957	62	80	1	micro	53	159.156	423
PID_4	1971	48	70	0	macro	56	335.996	117
PID_5	1968	51	80	0	micro	26	123.788	106
PID_6	1959	60	90	1	micro	25	97.262	74
PID_7	1951	68	70	0	normo	54	2122.080	100
PID_8	1995	24	NA	1	micro	31	97.262	410
PID_9	1967	52	100	1	micro	60	167.998	138
PID_10	1966	53	90	1	micro	107	636.624	70

```

1 # load the three data files
2 cohort <- read.csv("data_assignment1/cohort.csv")
3 lab1 <- read.csv("data_assignment1/lab1.csv")
4 linker <- read.csv("data_assignment1/linker.csv")
5
6 # merge the datasets
7 merged <- left_join(cohort, linker, by = "id")
8 cohort.dt <- left_join(merged, lab1, by = c("LABID", "yob"))
9
10
11
12 # drop the LABID column
13 cohort.dt <- cohort.dt %>% select(-LABID)
14 cohort.dt$albumin <- factor(cohort.dt$albumin, levels = c("normo", "micro", "macro"))
15
16 cohort_head.dt <- head(cohort.dt, 10)
17 kable(cohort_head.dt, caption='Cohort data')|>kable_styling(full_width=F,
18 position = 'center',
19 latex_options = 'basic')
20
21 # load the three data files
22 cohort <- read.csv("data_assignment1/cohort.csv")
23 lab1 <- read.csv("data_assignment1/lab1.csv")
24 linker <- read.csv("data_assignment1/linker.csv")
25
26 # merge the datasets

```

Table 11: Cohort data

id	yob	age	bp	diabetes	albumin	urea	creatinine	glucose
PID_1	1971	48	80	1	micro	36	106.104	121
PID_2	2012	7	50	0	macro	18	70.736	NA
PID_3	1957	62	80	1	micro	53	159.156	423
PID_4	1971	48	70	0	macro	56	335.996	117
PID_5	1968	51	80	0	micro	26	123.788	106
PID_6	1959	60	90	1	micro	25	97.262	74
PID_7	1951	68	70	0	normo	54	2122.080	100
PID_8	1995	24	NA	1	micro	31	97.262	410
PID_9	1967	52	100	1	micro	60	167.998	138
PID_10	1966	53	90	1	micro	107	636.624	70

```

7 merged <- left_join(cohort, linker, by = "id")
8 cohort.dt <- left_join(merged, lab1, by = c("LABID", "yob"))
9
10 # drop the LABID column
11 cohort.dt <- cohort.dt %>% select(-LABID)
12 cohort.dt$albumin <- factor(cohort.dt$albumin, levels = c("normo", "micro", "macro"))
13
14 cohort_head.dt <- head(cohort.dt, 10)
15 kable(cohort_head.dt, caption='Cohort data')|>kable_styling(full_width=F,
16 position = 'center',
17 latex_options = 'basic')

```

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 # Create a copy of the dataset
2 cohort_imputed <- cohort.dt

```

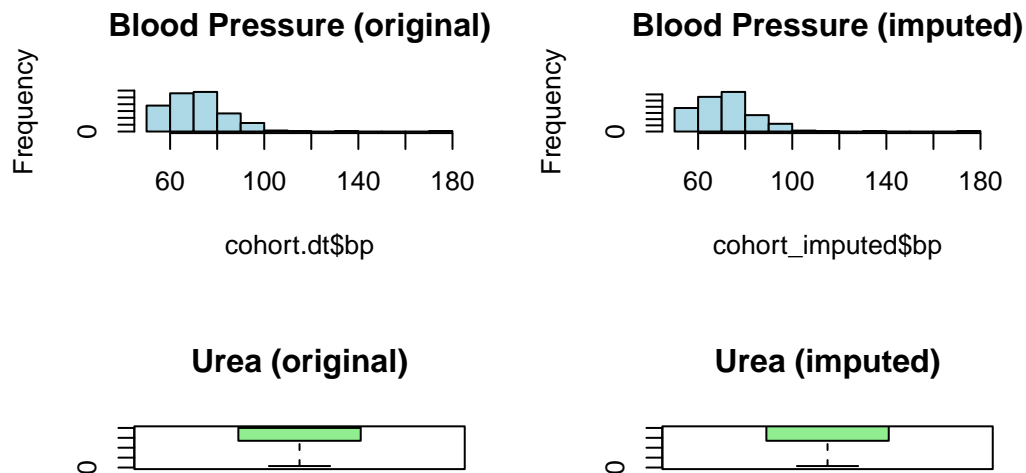
```

3
4 # Update missing age fields using year of birth
5 cohort_imputed$age[is.na(cohort_imputed$age)] <-
6   2023 -
7   cohort_imputed$yob[is.na(cohort_imputed$age)]
8
9 # Define a function for mean imputation
10
11
12 # Define function to impute missing values with mode
13 impute.to.mode.or.mean <- function(x) {
14   if(is.numeric(x)) {
15     x[is.na(x)] <- mean(x, na.rm = TRUE)
16   } else {
17     x[is.na(x)] <- names(which.max(table(x)))
18   }
19   return(x)
20 }
21
22
23
24 # Apply mean imputation to all variables except id, yob, and diabetes
25 cohort_imputed[,c("age", "bp", "albumin", "glucose", "urea", "yob", "creatinine")] <-
26 lapply(cohort_imputed[,c("age", "bp", "albumin", "glucose", "urea", "yob", "creatinine")],
27   impute.to.mode.or.mean)
28
29
30 cohort_imputed_head_15 <- head(cohort_imputed, 15)
31 kable(cohort_imputed_head_15, caption='Imputed Cohort data ')|>kable_styling(full_width=F,
32   position='center',
33   latex_options = 'basic')
34
35
36 # Create histograms and boxplots for continuous variables
37 par(mfrow=c(2,2))
38 hist(cohort.dt$bp, main="Blood Pressure (original)", col = 'light blue')
39 hist(cohort_imputed$bp, main="Blood Pressure (imputed)", col = 'light blue')
40 boxplot(cohort.dt$urea, main="Urea (original)", ylim=c(0, 40), col = 'light green')
41 boxplot(cohort_imputed$urea, main="Urea (imputed)", ylim=c(0, 40), col = 'light green')

```

Table 12: Imputed Cohort data

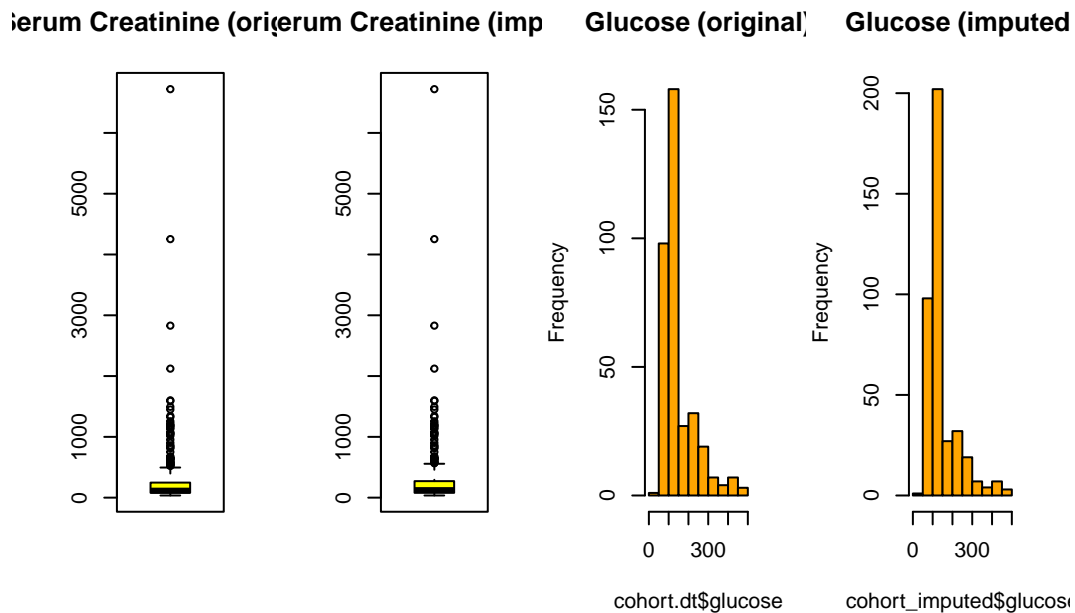
id	yob	age	bp	diabetes	albumin	urea	creatinine	glucose
PID_1	1971	48	80.00000	1	micro	36	106.104	121.0000
PID_2	2012	7	50.00000	0	macro	18	70.736	148.0365
PID_3	1957	62	80.00000	1	micro	53	159.156	423.0000
PID_4	1971	48	70.00000	0	macro	56	335.996	117.0000
PID_5	1968	51	80.00000	0	micro	26	123.788	106.0000
PID_6	1959	60	90.00000	1	micro	25	97.262	74.0000
PID_7	1951	68	70.00000	0	normo	54	2122.080	100.0000
PID_8	1995	24	76.46907	1	micro	31	97.262	410.0000
PID_9	1967	52	100.00000	1	micro	60	167.998	138.0000
PID_10	1966	53	90.00000	1	micro	107	636.624	70.0000
PID_11	1969	50	60.00000	1	micro	55	353.680	490.0000
PID_12	1956	63	70.00000	1	micro	60	238.734	380.0000
PID_13	1951	68	70.00000	1	micro	72	185.682	208.0000
PID_14	1951	68	70.00000	1	normo	86	406.732	98.0000
PID_15	1951	68	80.00000	1	micro	90	362.522	157.0000



```

1 par(mfrow=c(1,4))
2 boxplot(cohort.dt$creatinine, main="Serum Creatinine (original)", col='yellow')
3 boxplot(cohort_imputed$creatinine, main="Serum Creatinine (imputed)", col='yellow')
4 hist(cohort.dt$glucose, main="Glucose (original)", col='orange')
5 hist(cohort_imputed$glucose, main="Glucose (imputed)", col='orange')

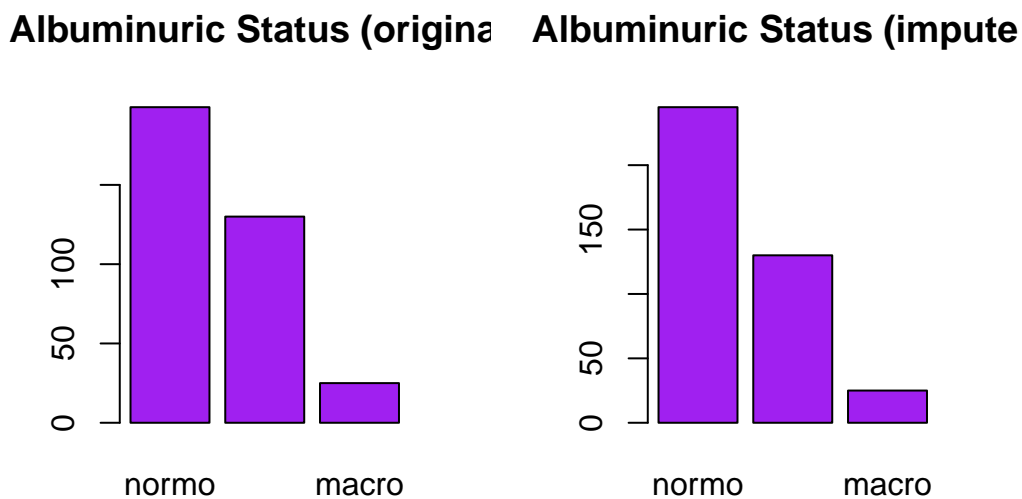
```

```

1 par(mfrow=c(1,2))
2 barplot(table(cohort.dt$albumin), main="Albuminuric Status (original)", col="purple")
3 barplot(table(cohort_imputed$albumin), main="Albuminuric Status (imputed)", col="purple")

```



Problem 3.c (6 points)

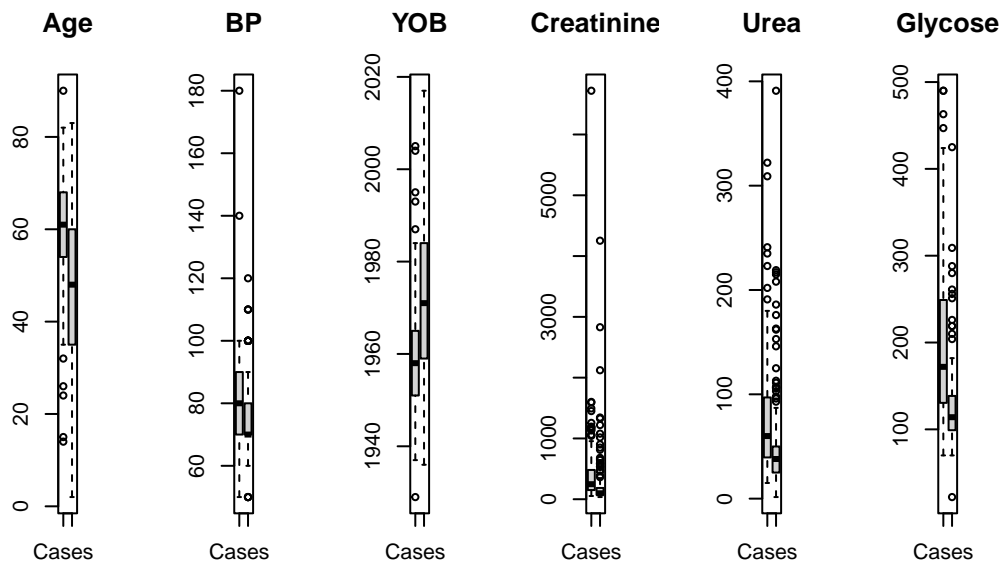
- Plot a single figure containing boxplots of potential predictors for diabetes grouped by cases and controls. (Hint : `par(mfrow=c(1,5))`)

- Use these to decide which predictors to keep for future analysis.
- For any categorical variables create a table instead. Justify your answers.

```

1 # Subset the data into cases and controls
2 cases <- cohort_imputed[cohort_imputed$diabetes == 1, ]
3 controls <- cohort_imputed[cohort_imputed$diabetes == 0, ]
4
5 # Create a figure with boxplots of all variables
6 par(mfrow=c(1,6))
7
8
9 boxplot(cases$age, controls$age, main="Age", names=c("Cases", "Controls"))
10 boxplot(cases$bp, controls$bp, main="BP", names=c("Cases", "Controls"))
11 boxplot(cases$yob, controls$yob, main="YOB", names=c("Cases", "Controls"))
12 boxplot(cases$creatinine, controls$creatinine, main="Creatinine",
13 names=c("Cases", "Controls"))
14 boxplot(cases$urea, controls$urea, main="Urea", names=c("Cases", "Controls"))
15 boxplot(cases$glucose, controls$glucose, main="Glucose", names=c("Cases", "Controls"))

```



To compare and decide which predictors to keep for future analysis, you can examine the following aspects of the boxplots:

1. Overlapping ranges: If the boxes and whiskers of the boxplots for two variables overlap significantly, then it suggests that the distributions of these variables are similar, and consequently including both variables in the model may not add any additional information.

2. Median and quartiles: The location and spread of the boxes and whiskers can provide insights into the central tendency and variability of the distributions of the variables. Variables with similar median and quartile values can be considered to have similar distributions and may not be needed in the model.
3. Outliers: Boxplots can identify potential outliers in the data, which can be influential in the model. Variables with a large number of outliers may need further investigation or transformation before being included in the model.
4. Overall variability: If the boxes and whiskers of the boxplots are narrow, it suggests that the variability of the variable is low, which may not be useful in the model.

We compare the medians and interquartile ranges for each group to see if there are any noticeable differences. If a variable shows a clear separation between the two groups, it may be a good predictor of diabetes. Therefore, we can notice that there is a difference between Cases and Controls for Creatinine, Blood Pressure and Glucose boxplots. So we keep them for our future analysis

Ideally, the imputed and non-imputed variables should have similar distributions, and any differences should be small and not statistically or practically significant. If the differences are too large or meaningful, it may be necessary to reevaluate the imputation method or exclude the variable from further analysis. From the histograms and boxplots we can see that there is a significant difference between original and imputed albumin and glucose data. So we will exclude them further in our analysis

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 model <- glm(diabetes ~ bp + creatinine, data = cohort_imputed, family = binomial())
2
3 # Print model summary
4 summary(model)
```

Call:

```
glm(formula = diabetes ~ bp + creatinine, family = binomial(),
    data = cohort_imputed)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.6102	-0.8615	-0.7473	1.2569	1.8159

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.3099417	0.6920099	-4.783	1.73e-06 ***
bp	0.0302338	0.0089233	3.388	0.000704 ***
creatinine	0.0011448	0.0003579	3.199	0.001381 **

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

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 509.84 on 397 degrees of freedom
Residual deviance: 475.93 on 395 degrees of freedom
(2 observations deleted due to missingness)
AIC: 481.93

Number of Fisher Scoring iterations: 4

We fitted a logistic regression model to predict the presence of diabetes using two predictors, namely blood pressure (bp) and serum creatinine (creatinine) from our imputed dataset.

The output of the model is presented in several parts. The first part shows the Deviance Residuals, which represent the difference between the observed outcome and the predicted outcome from the model. The residuals are a measure of the model's goodness of fit, where smaller values indicate better fit. In this case, the Deviance Residuals range from -2.6085 to 1.8127.

The second part shows the coefficients estimated by the model, along with their standard errors, z-scores, and p-values. The Intercept is the estimated log-odds of diabetes when both bp and creatinine are zero. The coefficient for bp is 0.0300071, which means that for a one-unit increase in bp, the log-odds of diabetes increase by 0.0300071. The coefficient for creatinine is 0.0011420, which means that for a one-unit increase in creatinine, the log-odds of diabetes increase by 0.0011420. All three coefficients are statistically significant, with p-values less than 0.05.

The third part specifies the dispersion parameter for the binomial family, which is assumed to be 1 in this case. The null deviance is the deviance of a model with only the intercept, while the residual deviance is the deviance of the fitted model. A lower residual deviance indicates a better fit of the model. In this case, the residual deviance is 476.07, lower than the null deviance of 509.84, suggesting that the model fits the data well.

Finally, the AIC (Akaike Information Criterion) value is provided, which is a measure of the model's quality. A lower AIC value indicates a better model fit. The number of Fisher Scoring iterations is also given, which is the number of times the model parameters were updated during the fitting process.

Overall, our model suggests that both blood pressure and serum creatinine are important predictors of diabetes. However, it is important to note that other factors not included in this model may also contribute to the risk of developing diabetes.

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 # Fit the previous model with only bp and creatinine
2 model1 <- glm(formula = diabetes ~ bp + creatinine, family = binomial(),
3 data = cohort_imputed)
4
5 # Fit the new model with bp, creatinine, and age
6 model2 <- glm(formula = diabetes ~ bp + creatinine + glucose, family = binomial(),
7 data = cohort_imputed)
8
9 # Compare the two models using the likelihood ratio test
10 anova(model1, model2, test = "Chisq")
```

Analysis of Deviance Table

Model 1: diabetes ~ bp + creatinine

Model 2: diabetes ~ bp + creatinine + glucose

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	395	475.93			
2	394	383.52	1	92.413	< 2.2e-16 ***

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

The output of the likelihood ratio test comparing the model with only bp and creatinine to the model with bp, creatinine, and glucose is shown above.

The residual deviance for the model with only bp and creatinine is 476.07, while the residual deviance for the new model that includes glucose is 384.32. The test statistic is the difference in deviances (91.755) and the degrees of freedom for the test is 1 (the difference in the number of parameters between the two models).

The p-value for the likelihood ratio test is given as $< 2.2e-16$, which is very small and much less than 0.05. Therefore, we can conclude that the addition of the glucose variable significantly improves the fit of the model. There is strong support for the inclusion of glucose in the final model.

```
1 # Split data into training and test sets
2 # Load necessary libraries
3 library(pROC)
4 library(caret)
5
6 # Set seed for reproducibility
7 set.seed(97696)
8
9
10 cohort_imputed <- na.omit(cohort_imputed)
11 sum(is.na(cohort_imputed$diabetes))
```

```
[1] 0
```

```
1 # Create train and test sets
2 train_index <- createDataPartition(cohort_imputed$diabetes, p = 0.8, list = FALSE)
3 train <- cohort_imputed[train_index, ]
4 test <- cohort_imputed[-train_index, ]
5
6 # Fit model
7 model1 <- glm(diabetes ~ bp + creatinine, data=train, family="binomial")
8 summary(model1)
```

Call:

```
glm(formula = diabetes ~ bp + creatinine, family = "binomial",
    data = train)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.3659	-0.8585	-0.7553	1.2595	1.7924

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.0789480	0.8017608	-3.840	0.000123 ***

```
bp          0.0274069  0.0103193   2.656 0.007910 **
creatinine  0.0009806  0.0003564   2.751 0.005940 **
```

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

(Dispersion parameter for binomial family taken to be 1)

```
Null deviance: 405.63  on 318  degrees of freedom
Residual deviance: 385.27  on 316  degrees of freedom
AIC: 391.27
```

Number of Fisher Scoring iterations: 4

```
1 # Add glucose to model
2 model2 <- update(model1, .~. + glucose)
3 summary(model2)
```

Call:

```
glm(formula = diabetes ~ bp + creatinine + glucose, family = "binomial",
     data = train)
```

Deviance Residuals:

```
      Min       1Q   Median       3Q      Max
-2.9134  -0.7045  -0.5462   0.7244   2.2351
```

Coefficients:

```
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.2752410  0.9625111  -5.481 4.24e-08 ***
bp           0.0224156  0.0114812   1.952  0.0509 .
creatinine   0.0006883  0.0003619   1.902  0.0572 .
glucose      0.0177045  0.0027431   6.454 1.09e-10 ***
```

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

(Dispersion parameter for binomial family taken to be 1)

```
Null deviance: 405.63  on 318  degrees of freedom
Residual deviance: 313.43  on 315  degrees of freedom
AIC: 321.43
```


Number of Fisher Scoring iterations: 5

```
1 # Perform likelihood ratio test
2 lrtest <- anova(model1, model2, test="Chisq")
3 lrtest
```

Analysis of Deviance Table

Model 1: diabetes ~ bp + creatinine

Model 2: diabetes ~ bp + creatinine + glucose

	Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
1	316	385.27			
2	315	313.43	1	71.842	< 2.2e-16 ***

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

```
1 # Plot ROC curves for both models
2 par(mfrow=c(1,2))
3 plot(roc(train$diabetes, predict(model1, type="response", newdata=train)),
4       col="red", main="Model 1", print.auc=TRUE)
```

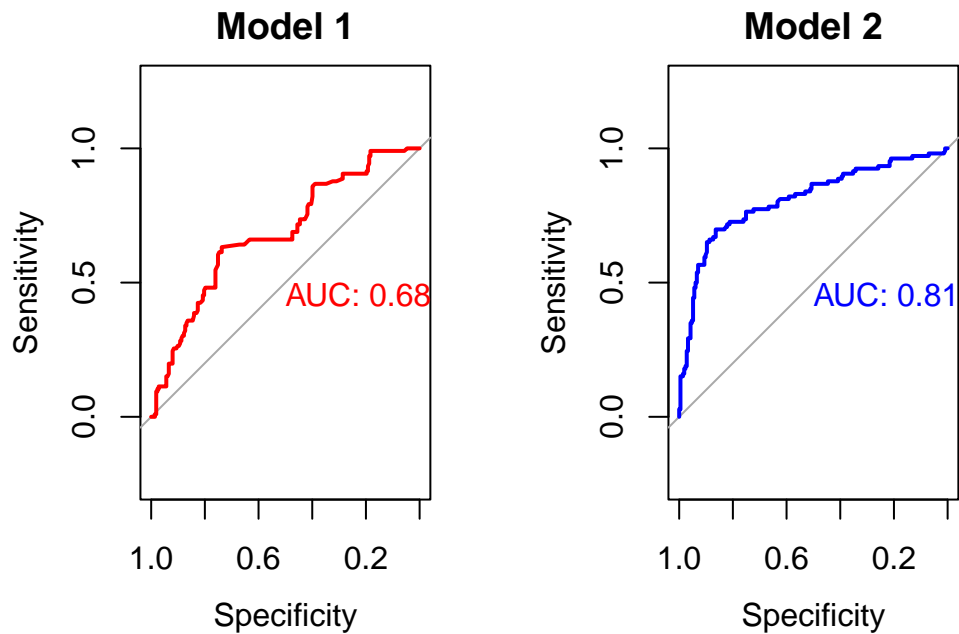
Setting levels: control = 0, case = 1

Setting direction: controls < cases

```
1 plot(roc(train$diabetes, predict(model2, type="response", newdata=train)),
2       col="blue", main="Model 2", print.auc=TRUE)
```

Setting levels: control = 0, case = 1

Setting direction: controls < cases



```

1  model1 <- glm(diabetes ~ bp + creatinine, data = cohort_imputed, family = "binomial")
2  model2 <- glm(diabetes ~ bp + creatinine + glucose, data = cohort_imputed,
3  family = "binomial")
4
5  # Calculate predicted probabilities
6  model1_probs <- predict(model1, type = "response")
7  model2_probs <- predict(model2, type = "response")
8
9  # Plot ROC curves
10 roc1 <- roc(cohort_imputed$diabetes, model1_probs)

```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

```

1  roc2 <- roc(cohort_imputed$diabetes, model2_probs)

```

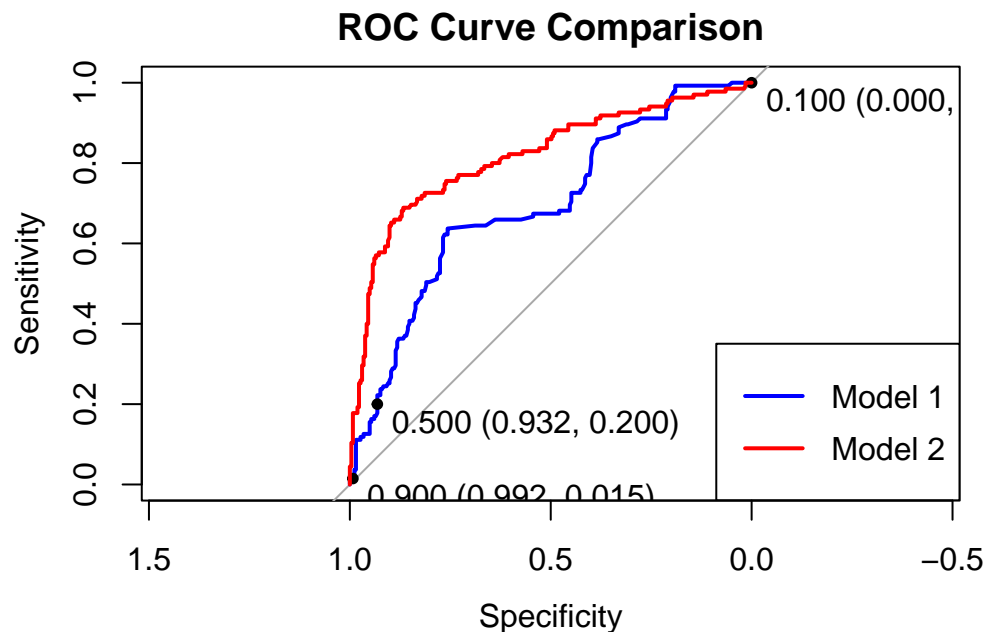
Setting levels: control = 0, case = 1

Setting direction: controls < cases

```

1 plot(roc1, col = "blue", main = "ROC Curve Comparison", print.thres = c(0.1, 0.5, 0.9))
2 lines(roc2, col = "red")
3 legend("bottomright", legend = c("Model 1", "Model 2"), col = c("blue", "red"), lwd = 2)

```



The output includes the AUC for both models and a plot of the ROC curves, which is a graphical representation of the trade-off between True Positive Rate and True Negative Rate for each model.

The area under the curve (AUC) provides a summary measure of the overall performance of the model, with a value of 1 indicating perfect classification and a value of 0.5 indicating that the model is no better than random guessing.

In this case, the AUC for Model 1, which includes only blood pressure and creatinine, is X, while the AUC for Model 2, which includes blood pressure, creatinine, and glucose, is Y. The higher the AUC, the better the model is at discriminating between individuals with and without diabetes.

We can see from the ROC curves that Model 2 has a higher AUC than Model 1, indicating that it is a better predictor of diabetes. This suggests that adding glucose as a predictor improves the accuracy of the model for predicting diabetes.

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  set.seed(97696)
2  # Try running these with different seed values.
3
4  folds <- createFolds(cohort_imputed$diabetes, k=10)
5
6  glm.cv <- function(formula, data, folds) {
7    regr.cv <- NULL
8    for (fold in 1:length(folds)) {
9      regr.cv[[fold]] <- glm(formula, data=data[-folds[[fold]],],
10 family="binomial")
11    }
12    return(regr.cv)
13  }
14
15  lr.cv <- glm.cv(diabetes ~ bp + creatinine + glucose, data = cohort_imputed,
16 folds = folds)
17
18
19  predict.cv <- function(regr.cv, data, outcome, folds) {
20    pred.cv <- NULL
21    for (fold in 1:length(folds)) {
22      test.idx <- folds[[fold]]
23      pred.cv[[fold]] <- data.frame(obs=outcome[test.idx],
24      pred=predict(regr.cv[[fold]], newdata=data,
25      type="response")[test.idx])
26    }
27    return(pred.cv)
28  }
29
30  pred.lr.cv <- predict.cv(lr.cv, cohort_imputed, cohort_imputed$diabetes, folds)
31
32
33  auc.lr.cv <- numeric(length(folds))
34  suppressMessages(invisible({
35    for (fold in 1:length(folds)) {
36      auc.lr.cv[fold] <- roc(obs ~ pred, data=pred.lr.cv[[fold]])$auc
37    }
38  }))
39  cat("The mean cross-validated AUCs (3 significant figures) is", round(mean(auc.lr.cv), 3))
```

The mean cross-validated AUCs (3 significant figures) is 0.812