# Statistics for Computer Science

## Assignment 1

## Martin Gregorík

## 500359

Services Development Management

Faculty of Informatics
Masaryk University

April 9, 2020

# Exercise 1

First, I loaded the data. Then I made a subset of cranial breadth of males from populations AUSTRALI and PERU. This was unsorted, so I sorted it and checked for missing values. There were no missing values.

```
1   # I commented setwd because you probably don't have this directory on
        your machine.
2   # And .Rnw needs to be compilable on your machine.
3   #setwd("/home/martingregorik/R/assignment01")
4   options(max.print=10000)
5   library(xtable)
6
7   howell <- read.csv("Howell.csv", header = TRUE)
8   #str(howell)
9
10  xcb_unsort <- howell$XCB[howell$Sex == 'M' & (howell$Population == '
        AUSTRALI' | howell$Population == 'PERU')]
11  xcb <- sort(xcb_unsort)
12  #sum(is.na(xcb)) = 0
13  #is.na(xcb)
```

Then I created custom functions.

```
14  # Task 1
15  my_sample_min <- function(vec) {
16     min <- Inf
17     for (i in 1:length(vec)) {
18       if (vec[i] < min) {
19         min <- vec[i]
20       }
21     }
22     return (min)
23  }
24
25  my_sample_max <- function(vec) {
26     max <- -Inf
27     for (i in 1:length(vec)) {
28       if (vec[i] > max) {
29         max <- vec[i]
30       }
31     }
32     return (max)
33  }
34
35  my_sample_mean <- function(vec) {
36     sum <- 0
37     for (i in 1:length(vec)) {
38       sum <- sum + vec[i]
39     }
```

```
40      return (sum / length(vec))
41  }
42
43  # sum of (xi - x~) squared
44  my_sum_sample_avg <- function(vec, exponent=2) {
45      sum <- 0
46      div <- 0
47      avg <- my_sample_mean(vec)
48      for (i in 1:length(vec)) {
49        div <- vec[i] - avg
50        sum <- sum + (div ^ exponent)
51      }
52      return (sum)
53  }
54
55  my_decile <- function(vec, k) {
56      # k / 10 * 100
57      return (vec[1:(k * 10)])
58  }
59
60  my_quartile <- function(vec, q) {
61      # denominator
62      denom <- 1 / q
63      len <- length(vec)
64      if (len %% 2 == 0) {
65        # even
66        return ((vec[len / denom] + vec[len / denom + 1]) / 2)
67      } else {
68        # odd
69        return (vec[(len + 1) / denom])
70      }
71  }
72  median <- my_quartile(xcb, 0.5)
73
74  my_five_num_sum <- function(vec) {
75      return (data.frame(min=my_sample_min(vec), lower_q=my_quartile(vec,
          0.25), median=my_quartile(vec, 0.50), upper_q=my_quartile(vec,
          0.75), max=my_sample_max(vec)))
76  }
77
78  my_sample_skew_cramer <- function(vec) {
79      nom <- my_sum_sample_avg(vec, 3)
80      denom <- length(vec) * (my_sample_variance(vec) ^ (3 / 2))
81      return (nom / denom)
82  }
83
84  my_sample_kurtosis <- function(vec) {
85      nom <- my_sum_sample_avg(vec, 4)
86      denom <- length(vec) * (my_sample_variance(vec) ^ 2)
```

```
 87    return ((nom / denom) - 3)
 88  }
 89  # broad = thick tails = platykurtic
 90
 91  my_sample_variance <- function(vec, exponent=2) {
 92    if (length(vec) != 0) {
 93      avg <- my_sample_mean(vec)
 94      sum <- my_sum_sample_avg(vec, exponent)
 95      return (sum / (length(vec) - 1))
 96    }
 97  }
 98
 99  my_sample_standard_deviation <- function(vec, exponent=2) {
100    if (length(vec) != 0) {
101      pw <- my_sample_variance(vec, exponent)
102      return (sqrt(pw))
103    }
104  }
105
106  my_sample_range <- function(vec) {
107    return (my_sample_max(vec) - my_sample_min(vec))
108  }
109
110  my_sample_decile_range <- function(vec) {
111    return (my_quartile(vec, 0.90) - my_quartile(vec, 0.10))
112  }
113
114  my_sample_trimmed_avg <- function(vec) {
115    gamma <- 0.1
116    n <- length(vec)
117    g <- floor(gamma * n)
118    # xtg
119    return ((1 / (n - 2 * g)) * sum(vec[g + 1:n - g]))
120  }
121
122  my_sample_trimmed_var <- function(vec) {
123    gamma <- 0.1
124    n <- length(vec)
125    g <- floor(gamma * n)
126    xtg <- my_sample_trimmed_avg(vec)
127    # stg
128    return ((1 / (n - (2 * g) - 1)) * sum((vec[g + 1:n - g] - xtg)^2))
129  }
```

Then I calculated characteristics of each population and stored them in a table.

```
130  # Australia
131  xcb_australi_unsorted <- howell$XCB[howell$Sex == 'M' & howell$
         Population == 'AUSTRALI']
132  xcb_australi <- sort(xcb_australi_unsorted)
```

```
133  xcb_aus_tab <- round(data.frame(
134    size=length(xcb_australi),
135    mean=my_sample_mean(xcb_australi),
136    my_five_num_sum(xcb_australi),
137    skew=my_sample_skew_cramer(xcb_australi),
138    kurt=my_sample_kurtosis(xcb_australi),
139    variance=my_sample_variance(xcb_australi),
140    sd=my_sample_standard_deviation(xcb_australi),
141    range=my_sample_range(xcb_australi),
142    dec_range=my_sample_decile_range(xcb_australi),
143    trim_avg=my_sample_trimmed_avg(xcb_australi),
144    trim_var=my_sample_trimmed_var(xcb_australi)
145  ), 4)
146  # Peru
147  xcb_peru_unsorted <- howell$XCB[howell$Sex == 'M' & howell$Population ==
         'PERU']
148  xcb_peru <- sort(xcb_peru_unsorted)
149  xcb_peru_tab <- round(data.frame(
150    size=length(xcb_peru),
151    mean=my_sample_mean(xcb_peru),
152    my_five_num_sum(xcb_peru),
153    skew=my_sample_skew_cramer(xcb_peru),
154    kurt=my_sample_kurtosis(xcb_peru),
155    variance=my_sample_variance(xcb_peru),
156    sd=my_sample_standard_deviation(xcb_peru),
157    range=my_sample_range(xcb_peru),
158    dec_range=my_sample_decile_range(xcb_peru),
159    trim_avg=my_sample_trimmed_avg(xcb_peru),
160    trim_var=my_sample_trimmed_var(xcb_peru)
161  ), 4)
162  # Concat them to single data frame
163  xcb_tab <- xcb_aus_tab
164  xcb_tab[2, ] <- xcb_peru_tab
165  rownames(xcb_tab) <- c("AUSTRALI", "PERU")
166  # Since the table is too long for pdf, I split them into two.
167  # Yes, I concatenated them earlier, but that was rows.
168  # Now I cut them in half by columns.
169  xcb_tab_first_half <- xcb_tab[, 1:(length(xcb_tab) / 2)]
170  xcb_tab_second_half <- xcb_tab[, ((length(xcb_tab) / 2) + 1):length(xcb_
      tab)]
```

|          | size  | mean   | min    | lower$_q$ | median | upper$_q$ | max    |
|----------|-------|--------|--------|-----------|--------|-----------|--------|
| AUSTRALI | 52.00 | 131.94 | 124.00 | 128.00    | 131.00 | 134.00    | 144.00 |
| PERU     | 55.00 | 137.95 | 129.00 | 135.00    | 138.00 | 141.00    | 149.00 |

Table 1: Characteristics of maximal cranial breadth of AUSTRALI and PERU populations

|          | skew  | kurt  | variance | sd   | range | dec$_r$ange | trim$_a$vg |
|----------|-------|-------|----------|------|-------|-------------|------------|
| AUSTRALI | 0.64  | -0.34 | 26.06    | 5.10 | 20.00 | 14.00       | 163.36     |
| PERU     | -0.00 | 0.06  | 15.87    | 3.98 | 20.00 | 9.00        | 168.60     |

Table 2: Characteristics of maximal cranial breadth of AUSTRALI and PERU populations

Next, I created boxplots of maximum cranial breadth for each population. I set the width of boxes to be proportional to sample sizes.

```
171  # Different lengths(australi is 52, peru 55), need to add 0s to the end(
         neglected in output)
172  n <- max(length(xcb_australi), length(xcb_peru))
173  xcb_australi_prolonged <- xcb_australi
174  length(xcb_australi_prolonged) <- n
175  max_b <- data.frame(AUSTRALI=xcb_australi_prolonged, PERU=xcb_peru)

176  # Variable australi_peru_cols will be used in following exercises
177  australi_peru_cols = c("dodgerblue4", "indianred")
178  boxplot(
179    max_b,
180    width = c(length(xcb_australi), length(xcb_peru)),
181    notch = TRUE,
182    main = "Boxplot of maximal cranial breadth",
183    xlab = "Countries",
184    ylab = "Maximal cranial breadth (mm)",
185    col = australi_peru_cols,
186    pch = 16
187  )
188  points(1:2, xcb_tab$mean, col = "green", pch = 16)
189  legend('topleft', pch = 16, legend = c("Maximum", "Average"), col = c("
         black", "green"))
```
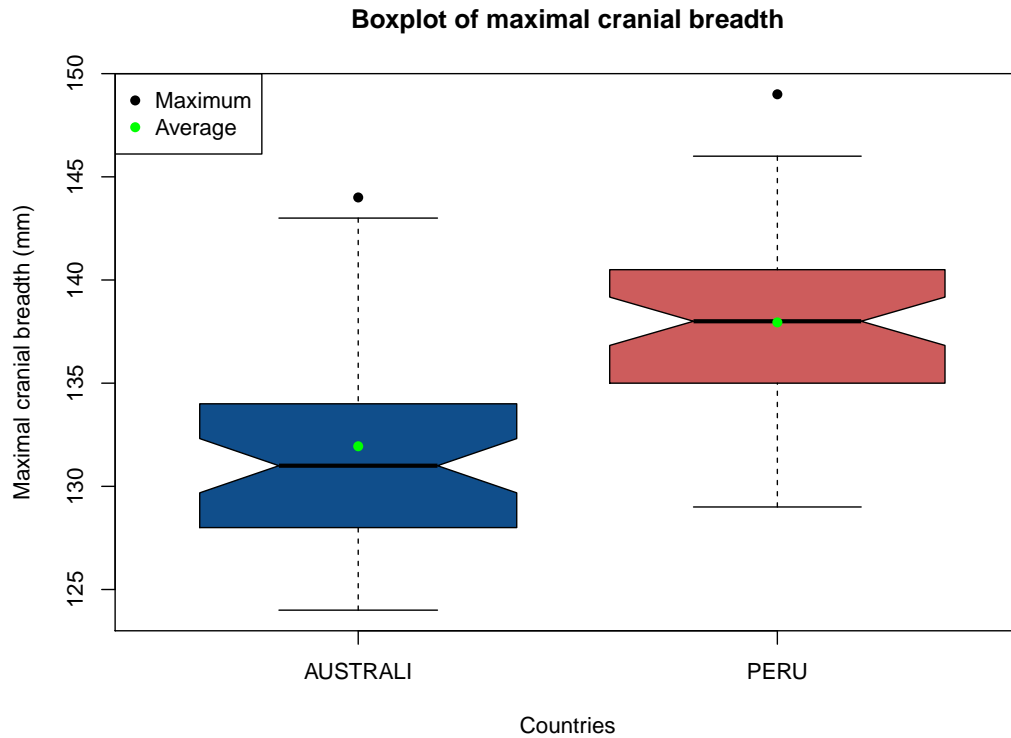
Figure 1: Boxplot of maximal cranial breadth(mm)

Created histogram of maximum cranial breadth for each population.

```
190  def.par <- par()
191  layout(matrix(c(1, 2), nrow=1, ncol=2), widths=c(0.5, 0.5))
192  #layout.show(n=2)
193  hist(
194    xcb_australi,
195    main = "Australi population",
196    xlab = "Maximal cranial breadth(mm)",
197    ylab = "Count",
198    col = australi_peru_cols[1]
199  )
200  hist(
201    xcb_peru,
202    main = "Peru population",
203    xlab = "Maximal cranial breadth(mm)",
204    ylab = "Count",
205    col = australi_peru_cols[2]
206  )
```
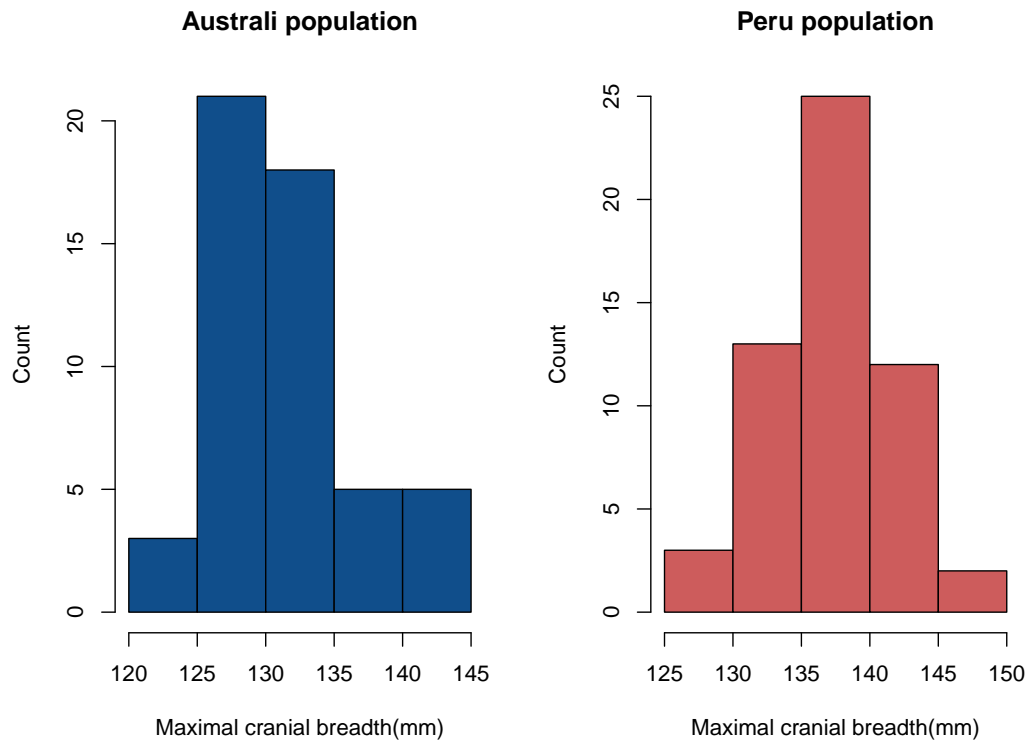
Figure 2: Histogram of maximal cranial breadth(mm)

Created normal qq-plot of maximum cranial breadth for each population.

The observed measurements(Y axis values) are from the same range so the plots can be easily compared.

```
207  layout(matrix(c(1, 2), nrow=1, ncol=2), widths=c(0.5, 0.5))
208  #layout.show(n=2)
209  qqnorm(
210    y = xcb_australi,
211    main = "Australi population",
212    xlab = "Theoretical quantities",
213    ylab = "Maximal cranial breadth(mm)",
214    ylim = c(min(xcb_australi, xcb_peru), max(xcb_australi, xcb_peru)),
215    pch = 16,
216    col = australi_peru_cols[1]
217  )
218  qqline(xcb_australi)
219  qqnorm(
220    y = xcb_peru,
221    main = "Peru population",
222    xlab = "Theoretical quantities",
223    ylab = "Maximal cranial breadth(mm)",
224    ylim = c(min(xcb_australi, xcb_peru), max(xcb_australi, xcb_peru)),
225    pch = 16,
226    col = australi_peru_cols[2]
```
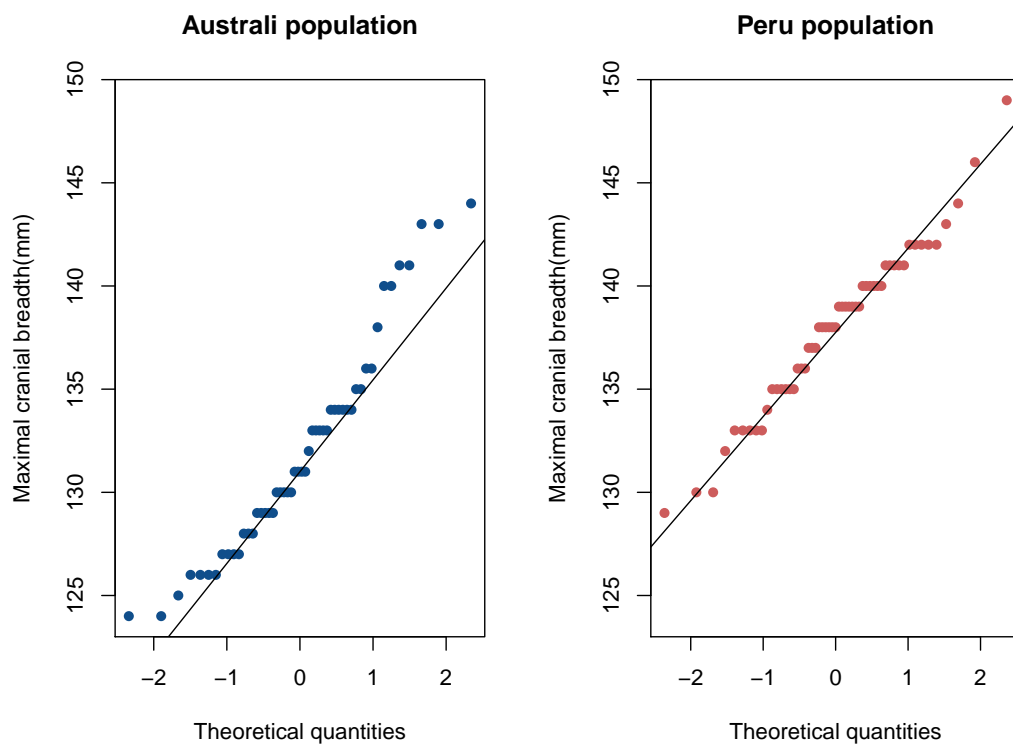
```
227  )
228  qqline ( xcb_peru )
```



Figure 3: Normal qq-plot of maximal cranial breadth(mm)

## Results and interpretation

Although both populations have same maximal cranial breadth reached, the differences can be best seen in histograms - more people of Peru population have bigger breadth than Australian people.

Variable XCB, representing maximal cranial breadth of each population, seems to be normally distributed variable in both cases.

# Exercise 2

First, I loaded the data. Then I checked for missing values in desired subset. There were none. So I summed the values in each column, which gave me a total sum of men and women in each year. Lastly, I added margins, which gave me total population in all of the years. Also, I like to keep variables at the beggining of each significant part, in this case Exercise.

```
229  area_esp <- read.csv("area_spanish_provinces.csv", header = TRUE)
230  pop_esp <- read.csv2("population-spain-1998-2018.csv", header = TRUE)
231  #str(area_esp)
232  #str(pop_esp)
233  #pop_esp
234  #sum(is.na(pop_esp[, 2:length(pop_esp)])) = 0
235  ### VARIABLES ###
236  legend_sex <- c("Women", "Men")
237  stat_years <-  c("1998", "2018")
238  total_pop_cols = c("pink", "dodgerblue")
239  year_cols = c("darkgoldenrod", "chocolate")
240  ### VARIABLES ###
241  people_each_year <- as.table(colSums(pop_esp[, 2:length(pop_esp)]))
242  total_people_each_year <- addmargins(people_each_year, FUN = c(Total=sum
       ))
```

Then I created a barplot of total population of Spain in each of the years, with each bar divided between men and women.
The data were stored as table(from previous task), so I converted them into data frame to be able to subset values.
The years were descending from 2018, so I made them ascending.
Then I stored the data in matrix in order for barplot to accept the argument type. Since the data in data frame is ordered by gender, I simply added "byrow = TRUE" to the matrix in order to store the data correctly.
I have also set a few plot parameters:
    "cex.axis", to make Y axis values smaller,
    "las", to make axis labels horizontal,
    "bty" in legend to hide legend box.

```
243  df_people <- as.data.frame(people_each_year)
244  # switched ordering of years
245  df_people <- df_people[dim(df_people)[1]:1, ]
246  df_people <- matrix(
247    df_people[, 2],
248    nrow = 2,
249    ncol = 5,
250    byrow = TRUE,
251    dimnames = list(c("F", "M"), c("1998", "2003", "2008", "2013", "2018")
         )
252  )
253  # It is clean now
254  suppressWarnings(par(def.par))
255  barplot(
```

```
256     height = df_people,
257     main = "Total population of Spain",
258     ylim = c(0, max(df_people)),
259     beside = TRUE,
260     las = 1,
261     cex.axis = 0.8,
262     col = total_pop_cols)
263 legend('topleft', pch = 15, legend = legend_sex, col = total_pop_cols,
        bty = 'n')
```
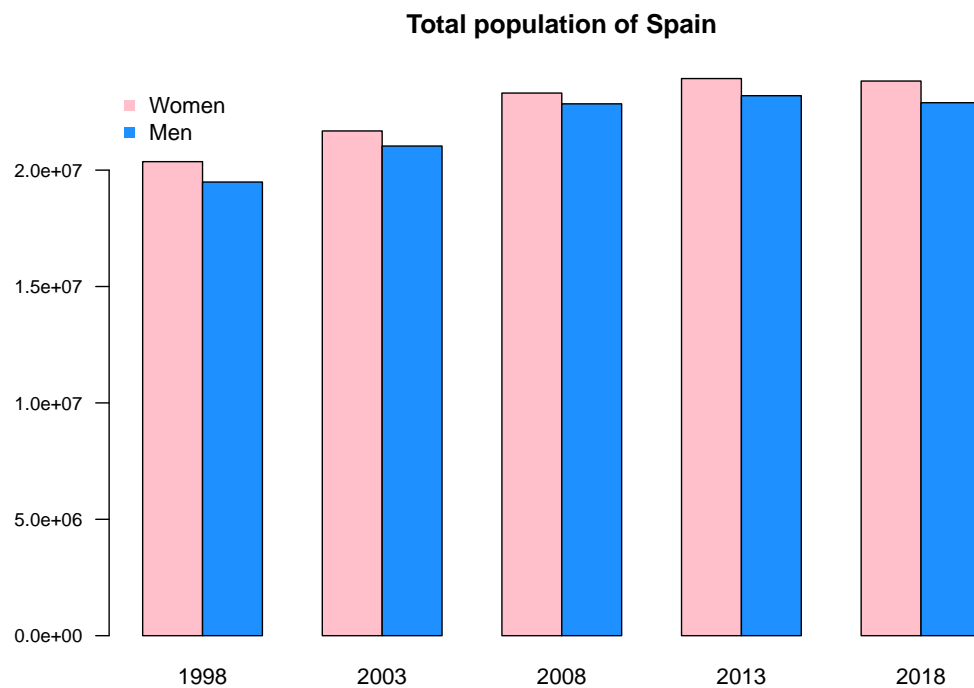


Figure 4: Barplot of total population in Spain

Then I created a barplot of relative proportions of men and women in each province in 2018. I transposed the input matrix in order to show the data by row(men and women side-by-side). I have also set few plot parameters:

"cex.names", to make the provinces' names smaller,

"las", to make axis labels perpendicular to the axis.

```
264  relative_pop_2018_provinces <- matrix(
265    data = c(pop_esp$females.2018, pop_esp$males.2018),
266    nrow = 52,
267    ncol = 2
268  )
269  barplot(
270    height = t(relative_pop_2018_provinces),
271    names.arg = pop_esp$province,
272    cex.names = 0.6,
273    cex.axis = 0.8,
274    main = "2018 relative population of Spain",
275    ylim = c(0, max(relative_pop_2018_provinces)),
276    las = 2,
277    beside = TRUE,
278    col = total_pop_cols)
279  legend('topleft', pch = 15, legend = legend_sex, col = total_pop_cols,
         bty = 'n')
```
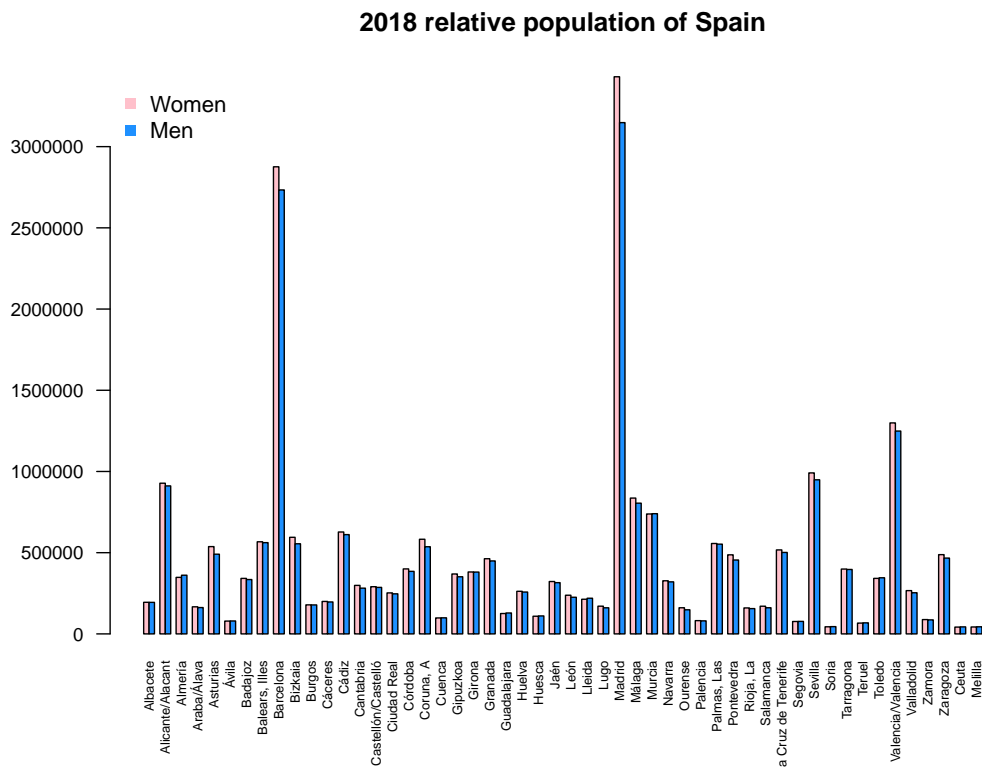


Figure 5: Barplot of relative population of Spain in 2018

Then I calculated population density in 1998 and 2018, for each province. The are of each province was loaded in the beginning of this exercise. I created a helper matrix to store intermediate values.

```
280  ppl_in_two_years <- matrix(
281    data = c(
282      pop_esp$males.1998,
283      pop_esp$females.1998,
284      pop_esp$males.2018,
285      pop_esp$females.2018
286      ),
287    nrow = nrow(pop_esp),
288    ncol = 4,
289    dimnames = list(pop_esp$province, c("males.1998", "females.1998", "
         males.2018", "females.2018"))
290  )
291  total_1998 <- as.vector(rowSums(ppl_in_two_years[, 1:2]))
292  total_2018 <- as.vector(rowSums(ppl_in_two_years[, 3:4]))
293  density_1998 <- total_1998 / area_esp$Area
294  density_2018 <- total_2018 / area_esp$Area
295  den <- data.frame(
296    Province=pop_esp$province,
297    "Year 1998" = density_1998,
298    "Year 2018" = density_2018
299  )
```

In subtask a, fistly, I sorted the data by year. Next, I calculated the estimates of characteristics of population density, each year separately. Lastly, I bound them together into one table.

```
300  # subtask a)
301  den_sorted_1998 <- sort(den$Year.1998)
302  den_sorted_1998_tab <- round(data.frame(
303    size=length(den_sorted_1998),
304    mean=my_sample_mean(den_sorted_1998),
305    my_five_num_sum(den_sorted_1998),
306    skew=my_sample_skew_cramer(den_sorted_1998),
307    kurt=my_sample_kurtosis(den_sorted_1998),
308    sd=my_sample_standard_deviation(den_sorted_1998)
309  ), 4)
310  den_sorted_2018 <- sort(den$Year.2018)
311  den_sorted_2018_tab <- round(data.frame(
312    size=length(den_sorted_2018),
313    mean=my_sample_mean(den_sorted_2018),
314    my_five_num_sum(den_sorted_2018),
315    skew=my_sample_skew_cramer(den_sorted_2018),
316    kurt=my_sample_kurtosis(den_sorted_2018),
317    sd=my_sample_standard_deviation(den_sorted_2018)
318  ), 4)
319  density_characteristics <- rbind(den_sorted_1998_tab, den_sorted_2018_
         tab)
320  row.names(density_characteristics) <- stat_years
```

| | size | mean | min | lower$_q$ | median | upper$_q$ | max | skew | kurt | sd |
|---|---|---|---|---|---|---|---|---|---|---|
| 1998 | 52.00 | 282.01 | 9.17 | 29.36 | 79.29 | 118.94 | 4623.69 | 4.12 | 16.24 | 848.58 |
| 2018 | 52.00 | 356.87 | 9.02 | 28.23 | 84.74 | 151.19 | 6644.92 | 4.38 | 19.20 | 1128.98 |

Table 3: Characteristics of population density of Spain in 1998 and 2018

In subtask b, I created a boxplot of population density in 1998 and in 2018. I have set the Y axis limits in order to see the graphs nicely, because there were some outliers which would squish the whole plot.

```
321  # subtask b)
322  suppressWarnings(par(def.par))
323  boxplot(
324    den[2:3],
325    notch = TRUE,
326    main = "Pop. density in 1998 & 2018",
327    ylab = "Population density",
328    ylim = c(0, 250),
329    pch = 16,
330    col = year_cols
331  )
```
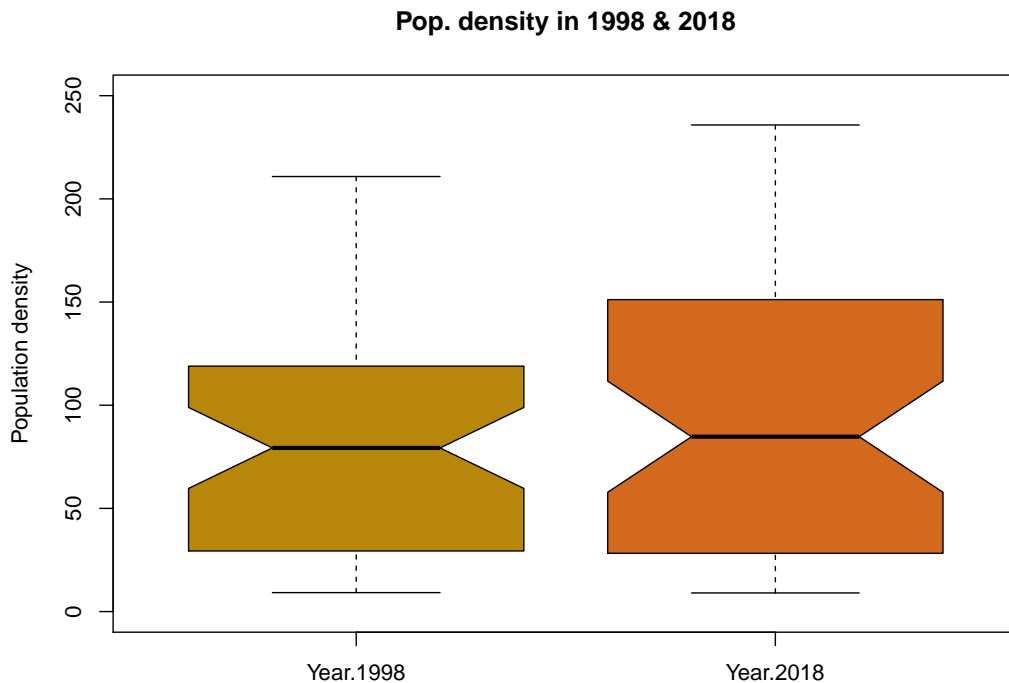


Figure 6: Boxplot of population density in Spain in 1998 and 2018

In subtask c, I created histogram of population density in 1998 and in 2018. There are two histograms, one for each year. Histogram for year 2018 is overlapping histogram for year 1998.

I changed the number of breaks. I was hoping to achieve a plot with province names, just like barplot of relative population of Spain(Figure 5), but it did not turn out that way. I'm really sad because that would make comparisons far better. Now I can't recognize the provinces. I left it as it is because I still think it looks better than the default.

Lastly, I also changed axes a little bit because I didn't like the values inserted by default, they were too sparse, too spread out.

```
332  # subtask c)
333  suppressWarnings(par(def.par))
334  rgb_cols = c(rgb(0.6,0.4,0.2,0.5), rgb(0.3,0.5,0.6,0.5))
335  hist(
336    den$Year.1998,
337    main = "Pop. density in 1998 & 2018",
338    xlab = "Density(people/km^2)",
339    ylab = "Number of provinces",
340    las = 2,
341    col = rgb_cols[1],
342    breaks = density_characteristics$size[1],
343    xaxt = 'n',
344    yaxt = 'n'
345  )
346  # X axis
347  axis(
348    side = 1,
349    at = seq(0, max(density_characteristics$max), by = 250),
350    las = 2,
351    cex.axis = 0.8
352  )
353  # Y axis
354  axis(
355    side = 2,
356    at = seq(0, density_characteristics$size[1], by = 5),
357    las = 2,
358    cex.axis = 0.8
359  )
360  hist(
361    den$Year.2018,
362    col = rgb_cols[2],
363    breaks = density_characteristics$size[2],
364    add = T)
365  legend('topright', pch = 15, legend = stat_years, col = rgb_cols, bty =
          'n')
```
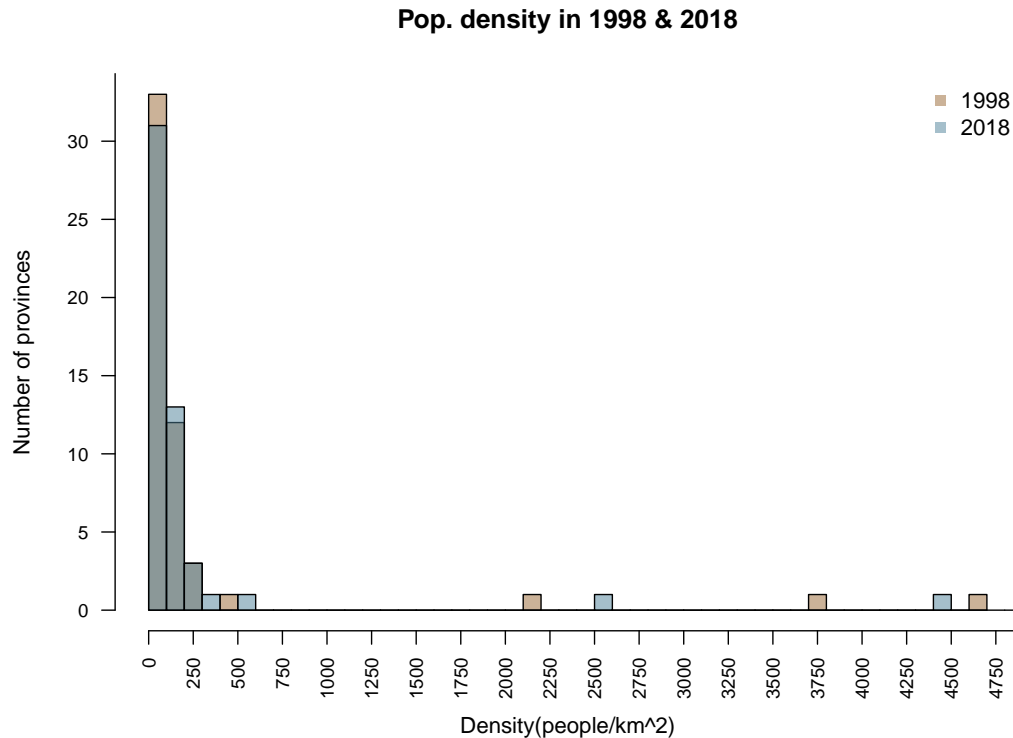
**Pop. density in 1998 & 2018**



Figure 7: Histogram of population density in Spain in 1998 and 2018

## Results and interpretation

Total population of Spain increased in years 1998-2018, even though there is a minor decrease between years 2013 and 2018.

Women are still slightly dominant when it comes to population proportions.

As we can see from the population density characteristics, the maximum population density increased from 4600 to 6600, which can also be seen in histogram. The mean population density increased from 282 to 356, which can be seen better e.g. in boxplot.