

Statistics for Computer Science

Assignment 1

Martin Gregorík

500359

Services Development Management

Faculty of Informatics
Masaryk University

April 10, 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)
14 ### VARIABLES ###
15 australi.peru.cols = c("dodgerblue4", "indianred")
16 ### VARIABLES ###

```

Then I created custom functions.

```

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

```

```
40   for (i in 1:length(vec)) {
41     sum <- sum + vec[i]
42   }
43   return (sum / length(vec))
44 }
45
46 # sum of (xi - x~) squared
47 MySumSampleAvg <- function(vec, exponent=2) {
48   sum <- 0
49   div <- 0
50   avg <- MySampleMean(vec)
51   for (i in 1:length(vec)) {
52     div <- vec[i] - avg
53     sum <- sum + (div ^ exponent)
54   }
55   return (sum)
56 }
57
58 MyDecile <- function(vec, k) {
59   # k / 10 * 100
60   return (vec[1:(k * 10)])
61 }
62
63 MyQuartile <- function(vec, q) {
64   # denominator
65   denom <- 1 / q
66   len <- length(vec)
67   if (len %% 2 == 0) {
68     # even
69     return ((vec[len / denom] + vec[len / denom + 1]) / 2)
70   } else {
71     # odd
72     return (vec[(len + 1) / denom])
73   }
74 }
75 median <- MyQuartile(xcb, 0.5)
76
77 MyFiveNumSum <- function(vec) {
78   return (data.frame(min=MySampleMin(vec), lower.q=MyQuartile(vec, 0.25)
79     , median=MyQuartile(vec, 0.50), upper.q=MyQuartile(vec, 0.75), max=
80     MySampleMax(vec)))
81 }
82
83 MySampleSkewCramer <- function(vec) {
84   nom <- MySumSampleAvg(vec, 3)
85   denom <- length(vec) * (MySampleVariance(vec) ^ (3 / 2))
86   return (nom / denom)
87 }
```

```

87 MySampleKurtosis <- function(vec) {
88   nom <- MySumSampleAvg(vec, 4)
89   denom <- length(vec) * (MySampleVariance(vec) ^ 2)
90   return ((nom / denom) - 3)
91 }
92 # broad = thick tails = platykurtic
93
94 MySampleVariance <- function(vec, exponent=2) {
95   if (length(vec) != 0) {
96     avg <- MySampleMean(vec)
97     sum <- MySumSampleAvg(vec, exponent)
98     return (sum / (length(vec) - 1))
99   }
100 }
101
102 MySampleSd <- function(vec, exponent=2) {
103   if (length(vec) != 0) {
104     pw <- MySampleVariance(vec, exponent)
105     return (sqrt(pw))
106   }
107 }
108
109 MySampleRange <- function(vec) {
110   return (MySampleMax(vec) - MySampleMin(vec))
111 }
112
113 MySampleDecileRange <- function(vec) {
114   return (MyQuartile(vec, 0.90) - MyQuartile(vec, 0.10))
115 }
116
117 MySampleTrimmedAvg <- function(vec) {
118   gamma <- 0.1
119   n <- length(vec)
120   g <- floor(gamma * n)
121   # xtg
122   return ((1 / (n - 2 * g)) * sum(vec[g + 1:n - g]))
123 }
124
125 MySampleTrimmedVar <- function(vec) {
126   gamma <- 0.1
127   n <- length(vec)
128   g <- floor(gamma * n)
129   xtg <- MySampleTrimmedAvg(vec)
130   # stg
131   return ((1 / (n - (2 * g) - 1)) * sum((vec[g + 1:n - g] - xtg)^2))
132 }

```

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

```

133 # Australia

```

```
134 xcb.australi.unsorted <- howell$XCB[howell$Sex == 'M' & howell$
    Population == 'AUSTRALI']
135 xcb.australi <- sort(xcb.australi.unsorted)
136 xcb.aus.tab <- round(data.frame(
137   size=length(xcb.australi),
138   mean=MySampleMean(xcb.australi),
139   MyFiveNumSum(xcb.australi),
140   skew=MySampleSkewCramer(xcb.australi),
141   kurt=MySampleKurtosis(xcb.australi),
142   variance=MySampleVariance(xcb.australi),
143   sd=MySampleSd(xcb.australi),
144   range=MySampleRange(xcb.australi),
145   dec.range=MySampleDecileRange(xcb.australi),
146   trim.avg=MySampleTrimmedAvg(xcb.australi),
147   trim.var=MySampleTrimmedVar(xcb.australi)
148 ), 4)
149 # Peru
150 xcb.peru.unsorted <- howell$XCB[howell$Sex == 'M' & howell$Population ==
    'PERU']
151 xcb.peru <- sort(xcb.peru.unsorted)
152 xcb.peru.tab <- round(data.frame(
153   size=length(xcb.peru),
154   mean=MySampleMean(xcb.peru),
155   MyFiveNumSum(xcb.peru),
156   skew=MySampleSkewCramer(xcb.peru),
157   kurt=MySampleKurtosis(xcb.peru),
158   variance=MySampleVariance(xcb.peru),
159   sd=MySampleSd(xcb.peru),
160   range=MySampleRange(xcb.peru),
161   dec.range=MySampleDecileRange(xcb.peru),
162   trim.avg=MySampleTrimmedAvg(xcb.peru),
163   trim.var=MySampleTrimmedVar(xcb.peru)
164 ), 4)
165 # Concat them to single data frame
166 xcb.tab <- xcb.aus.tab
167 xcb.tab[2, ] <- xcb.peru.tab
168 rownames(xcb.tab) <- c("AUSTRALI", "PERU")
169 # Since the table is too long for pdf, I split them into two.
170 # Yes, I concatenated them earlier, but that was rows.
171 # Now I cut them in half by columns.
172 xcb.tab.first.half <- xcb.tab[, 1:(length(xcb.tab) / 2)]
173 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.range	trim.avg
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.

```

174 # Different lengths (australi is 52, peru 55), need to add 0s to the end(
    neglected in output)
175 n <- max(length(xcb.australi), length(xcb.peru))
176 xcb.australi.prolonged <- xcb.australi
177 length(xcb.australi.prolonged) <- n
178 max.xcb <- data.frame(AUSTRALI=xcb.australi.prolonged, PERU=xcb.peru)

179 boxplot(
180   max.xcb,
181   width = c(length(xcb.australi), length(xcb.peru)),
182   notch = TRUE,
183   main = "Boxplot of maximal cranial breadth",
184   xlab = "Population",
185   ylab = "Maximal cranial breadth (mm)",
186   col = australi.peru.cols,
187   pch = 16
188 )
189 points(1:2, xcb.tab$mean, col = "green", pch = 16)
190 legend('topleft', pch = 16, legend = c("Maximum", "Average"), col = c("
    black", "green"))

```

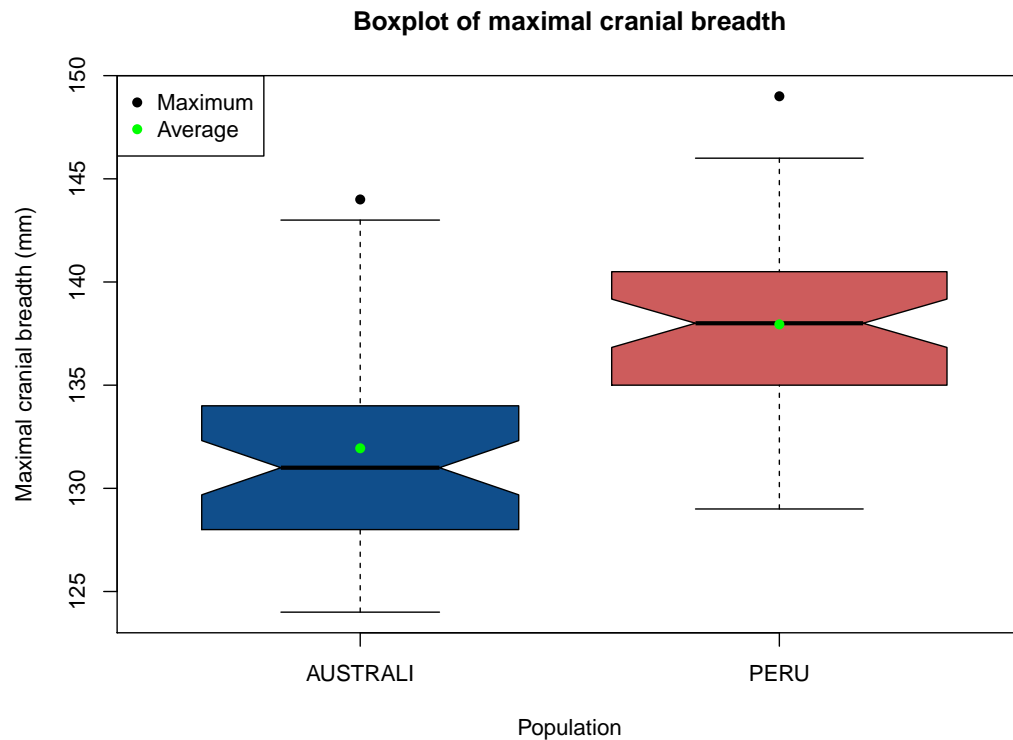


Figure 1: Boxplot of maximal cranial breadth(mm)

Created histogram of maximum cranial breadth for each population.

```

191 def.par <- par()
192 layout(matrix(c(1, 2), nrow=1, ncol=2), widths=c(0.5, 0.5))
193 #layout.show(n=2)
194 hist(
195   xcb.australi,
196   main = "Australi population",
197   xlab = "Maximal cranial breadth(mm)",
198   ylab = "Count",
199   col = australi.peru.cols[1]
200 )
201 hist(
202   xcb.peru,
203   main = "Peru population",
204   xlab = "Maximal cranial breadth(mm)",
205   ylab = "Count",
206   col = australi.peru.cols[2]
207 )

```

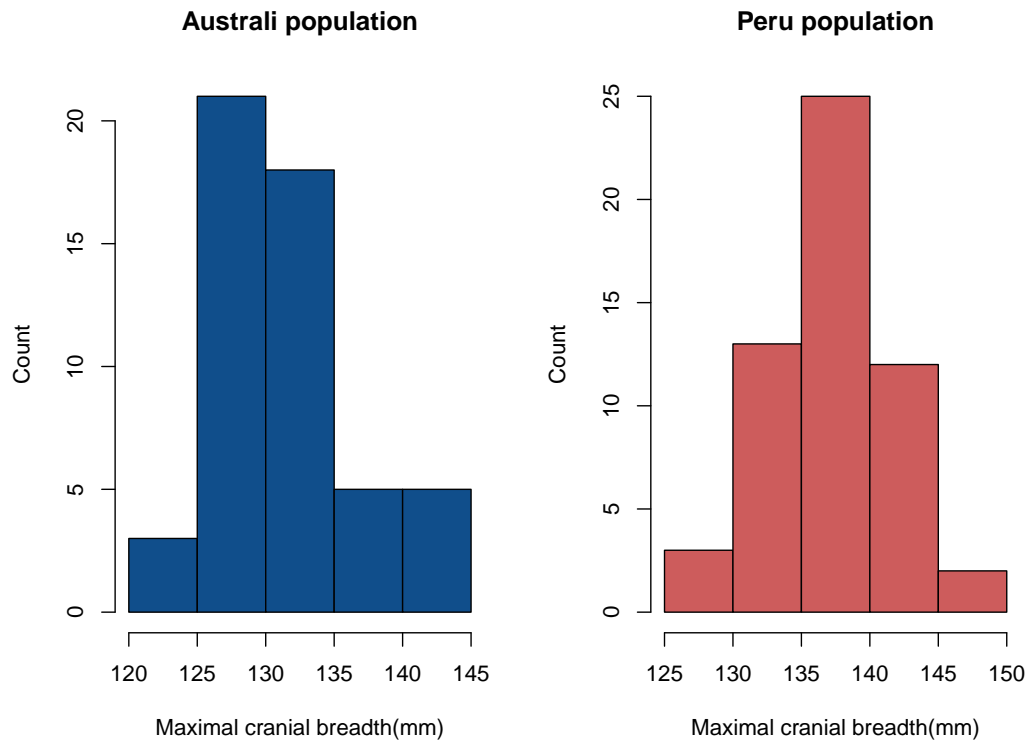


Figure 2: Histogram of maximal cranial breadth(mm)

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.

```

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

```



```
228 )  
229 qqline(xcb.peru)
```

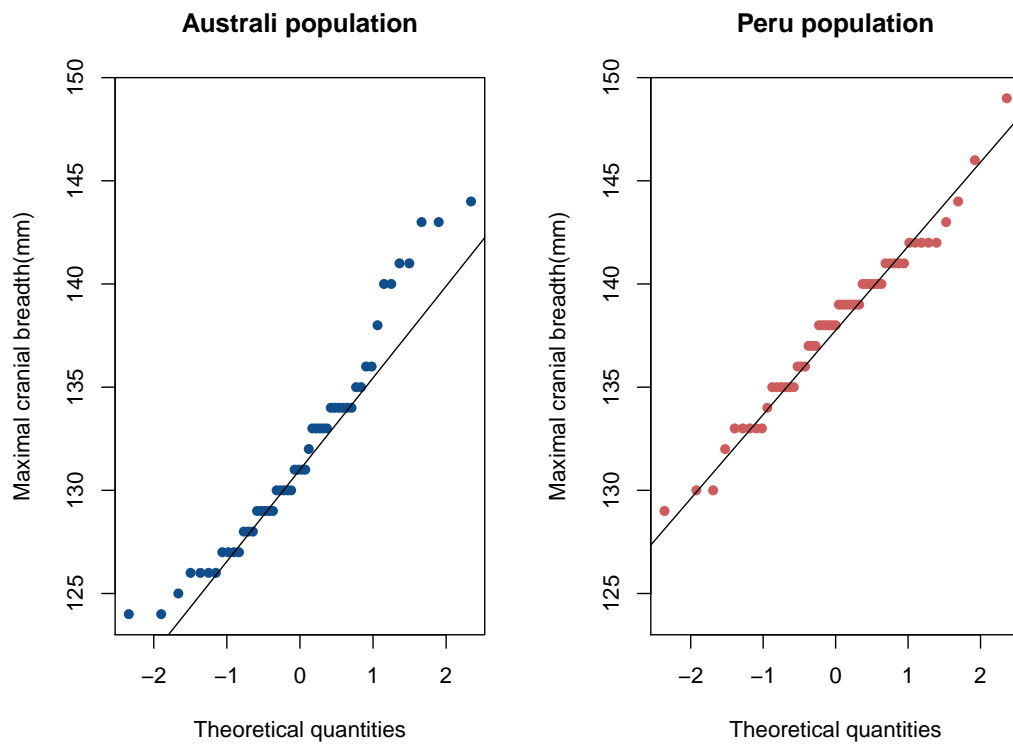


Figure 3: Normal qq-plot of maximal cranial breadth(mm)

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 beginning of each significant part, in this case Exercise.

```

230 area.esp <- read.csv("area_spanish_provinces.csv", header = TRUE)
231 pop.esp <- read.csv2("population-spain-1998-2018.csv", header = TRUE)
232 #str(area.esp)
233 #str(pop.esp)
234 #pop.esp
235 #sum(is.na(pop.esp[, 2:length(pop.esp)])) = 0
236 ### VARIABLES ###
237 legend.sex <- c("Women", "Men")
238 stat.years <- c("1998", "2018")
239 total.pop.cols = c("pink", "dodgerblue")
240 year.cols = c("darkgoldenrod", "chocolate")
241 ### VARIABLES ###
242 people.each.year <- as.table(colSums(pop.esp[, 2:length(pop.esp)]))
243 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.

```

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

```

```
257 height = df.people,
258 main = "Total population of Spain",
259 ylim = c(0, max(df.people)),
260 beside = TRUE,
261 las = 1,
262 cex.axis = 0.8,
263 col = total.pop.cols
264 )
265 legend('topleft', pch = 15, legend = legend.sex, col = total.pop.cols,
      bty = 'n')
```

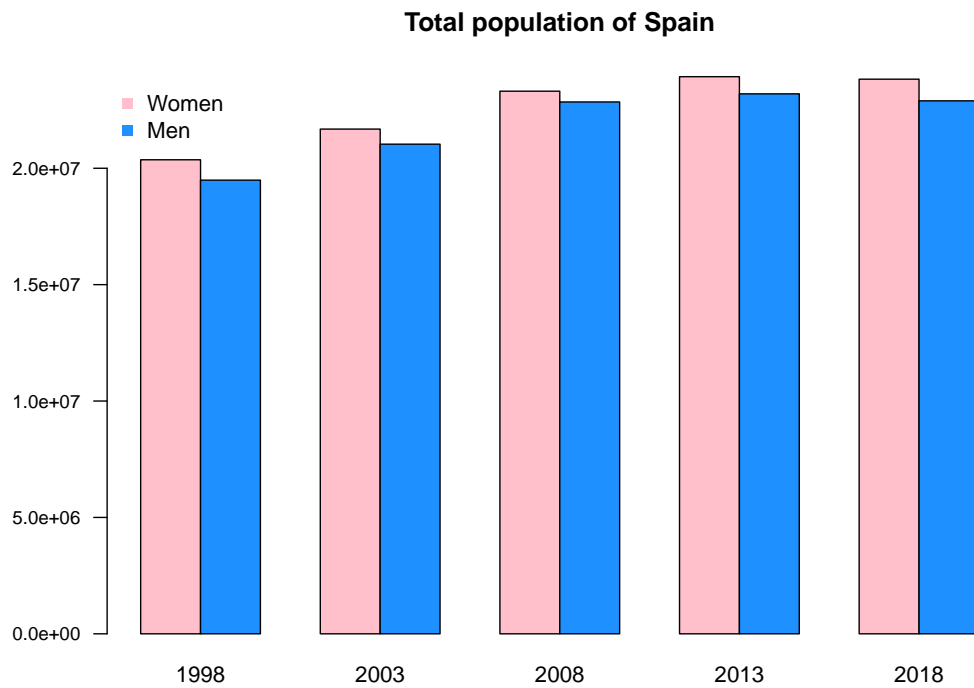


Figure 4: Barplot of total population in Spain

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.

```
266 relative.pop.2018.provinces <- matrix(  
267   data = c(pop.esp$females.2018, pop.esp$males.2018),  
268   nrow = 52,  
269   ncol = 2  
270 )  
271 barplot(  
272   height = t(relative.pop.2018.provinces),  
273   names.arg = pop.esp$province,  
274   cex.names = 0.6,  
275   cex.axis = 0.8,  
276   main = "2018 relative population of Spain",  
277   ylim = c(0, max(relative.pop.2018.provinces)),  
278   las = 2,  
279   beside = TRUE,  
280   col = total.pop.cols  
281 )  
282 legend('topleft', pch = 15, legend = legend.sex, col = total.pop.cols,  
      bty = 'n')
```

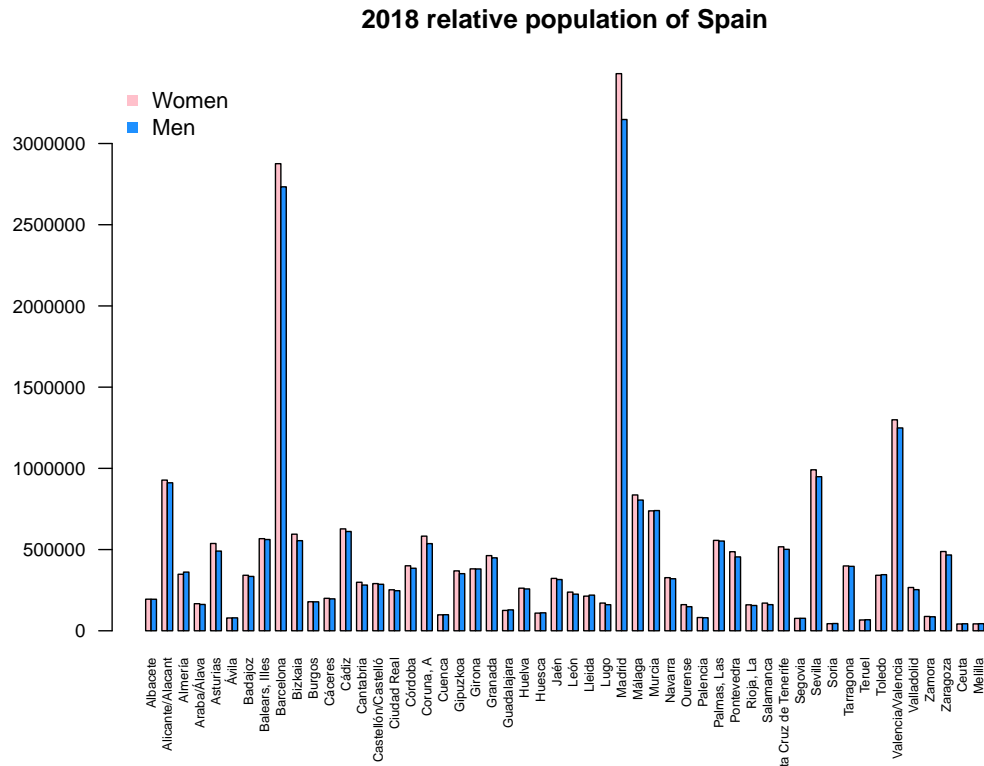


Figure 5: Barplot of relative population of Spain in 2018

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

```

283 ppl.in.two.years <- matrix(
284   data = c(
285     pop.esp$males.1998,
286     pop.esp$females.1998,
287     pop.esp$males.2018,
288     pop.esp$females.2018
289   ),
290   nrow = nrow(pop.esp),
291   ncol = 4,
292   dimnames = list(pop.esp$province, c("males.1998", "females.1998", "
      males.2018", "females.2018")))
293 )
294 total.1998 <- as.vector(rowSums(ppl.in.two.years[, 1:2]))
295 total.2018 <- as.vector(rowSums(ppl.in.two.years[, 3:4]))
296 density.1998 <- total.1998 / area.esp$Area
297 density.2018 <- total.2018 / area.esp$Area
298 den <- data.frame(
299   Province=pop.esp$province,
300   "Year 1998" = density.1998,
301   "Year 2018" = density.2018

```

302)

In subtask a, firstly, 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.

```

303 # subtask a)
304 den.sorted.1998 <- sort(den$Year.1998)
305 den.sorted.1998.tab <- round(data.frame(
306   size=length(den.sorted.1998),
307   mean=MySampleMean(den.sorted.1998),
308   MyFiveNumSum(den.sorted.1998),
309   skew=MySampleSkewCramer(den.sorted.1998),
310   kurt=MySampleKurtosis(den.sorted.1998),
311   sd=MySampleSd(den.sorted.1998)
312 ), 4)
313 den.sorted.2018 <- sort(den$Year.2018)
314 den.sorted.2018.tab <- round(data.frame(
315   size=length(den.sorted.2018),
316   mean=MySampleMean(den.sorted.2018),
317   MyFiveNumSum(den.sorted.2018),
318   skew=MySampleSkewCramer(den.sorted.2018),
319   kurt=MySampleKurtosis(den.sorted.2018),
320   sd=MySampleSd(den.sorted.2018)
321 ), 4)
322 density.characteristics <- rbind(den.sorted.1998.tab, den.sorted.2018.
   tab)
323 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.

```

324 # subtask b)
325 suppressWarnings(par(def.par))
326 boxplot(
327   den[2:3],
328   notch = TRUE,
329   main = "Pop. density in 1998 & 2018",
330   ylab = "Population density",
331   ylim = c(0, 250),
332   pch = 16,
333   col = year.cols
334 )

```

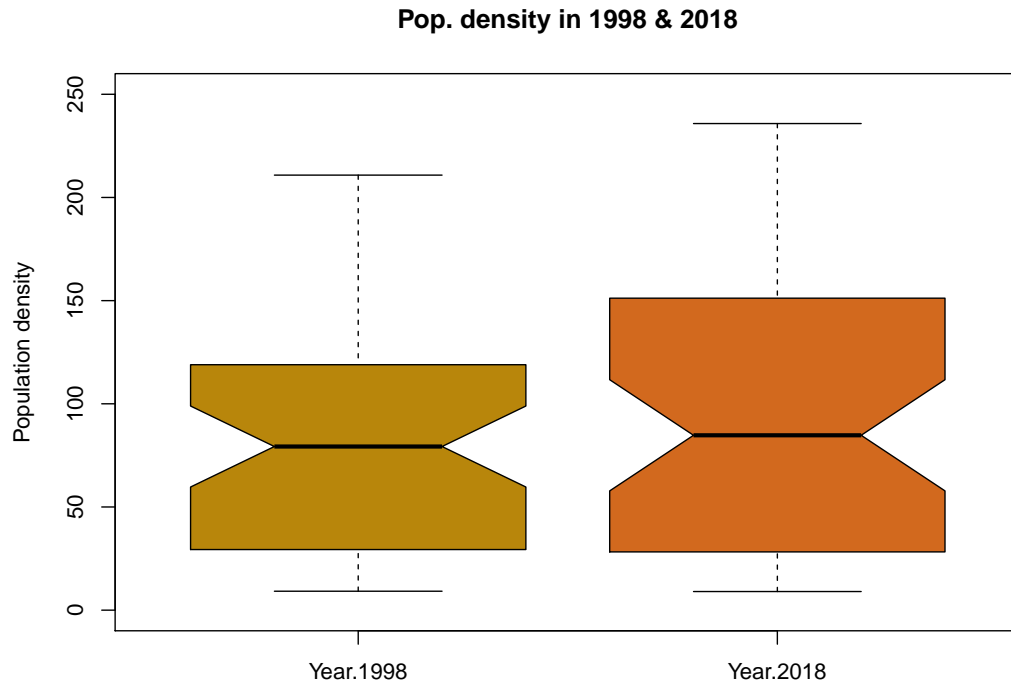


Figure 6: Boxplot of population density in Spain in 1998 and 2018

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.

```

335 # subtask c)
336 suppressWarnings(par(def.par))
337 rgb.cols = c(rgb(0.6,0.4,0.2,0.5), rgb(0.3,0.5,0.6,0.5))
338 hist(
339   den$Year.1998,
340   main = "Pop. density in 1998 & 2018",
341   xlab = "Density(people/km^2)",
342   ylab = "Number of provinces",
343   las = 2,
344   col = rgb.cols[1],
345   breaks = density.characteristics$size[1],
346   xaxt = 'n',
347   yaxt = 'n'
348 )
349 # X axis

```



```

350 axis(
351   side = 1,
352   at = seq(0, max(density.characteristics$max), by = 250),
353   las = 2,
354   cex.axis = 0.8
355 )
356 # Y axis
357 axis(
358   side = 2,
359   at = seq(0, density.characteristics$size[1], by = 5),
360   las = 2,
361   cex.axis = 0.8
362 )
363 hist(
364   den$Year.2018,
365   col = rgb.cols[2],
366   breaks = density.characteristics$size[2],
367   add = T)
368 legend('topright', pch = 15, legend = stat.years, col = rgb.cols, bty =
      'n')

```

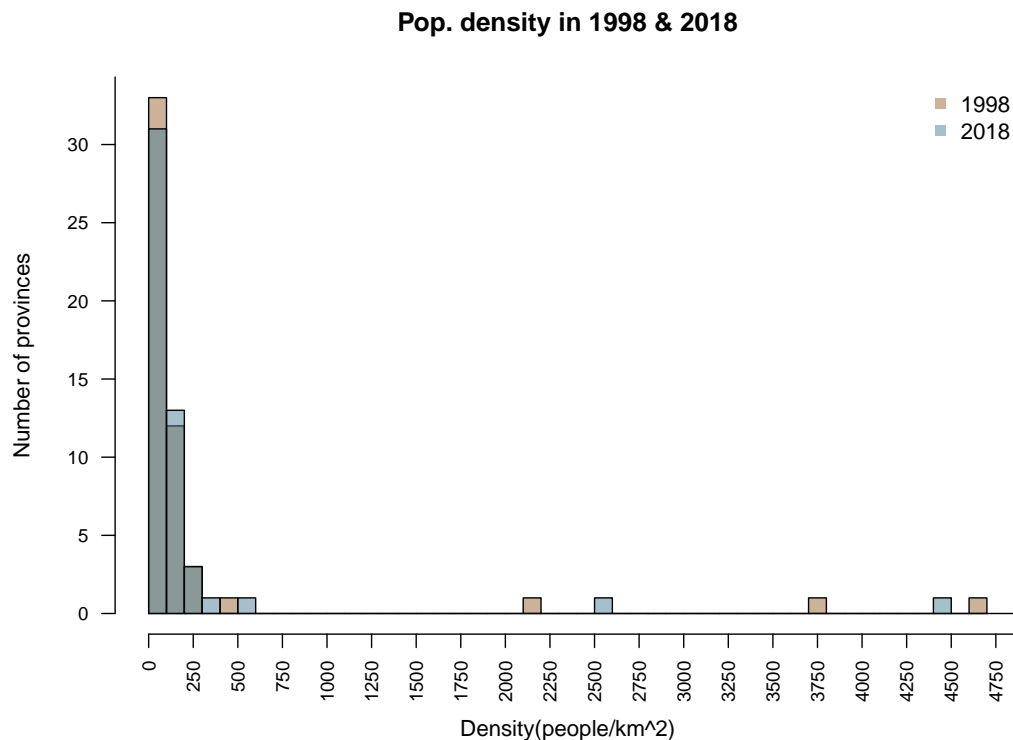


Figure 7: Histogram of population density in Spain in 1998 and 2018

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.