# Statistics for Computer Science

## Assignment 1

Kanitha Chim

501453

Field of Study Software System and Service Management

Faculty of Informatics
Masaryk University

April 10, 2020

# Exercise 1

1. Create your own R-functions for calculating estimates of the following characteristics: sample mean, sample five number summary, sample skewness, sample kurtosis, sample variance, sample standard deviation, sample range, sample decile range, sample 0.1-trimmed average and sample 0.1-trimmed variance. Use Cramér's estimates for skewness and kurtosis. Do not use inbuilt functions min(), max(), mean(), quantile(), var(), sd(), range() or functions from external libraries.

## Implementation in R

```
1  #function to calculate minimum
2  cal_minimum <- function(x){
3    x <- sort(x)
4    min_val <- x[1]
5    return(min_val)
6  }
7
8  #function to calculate maximum
9  cal_maximum <- function(x){
10    x <- sort(x)
11    max_val <- tail(x,1)
12    return(max_val)
13  }
14
15  #median, first quartile, thrid quartile, quartile_type should be, one,
      two(median), or three
16  cal_quartile <- function(x, quartile_type){
17    x <- sort(x)
18    result <- 0
19
20    if(quartile_type == 'two'){
21      if(length(x) %% 2 == 0){
22        return((x[length(x)/2] + x[length(x+1)/2])/2)
23      }else{
24        return(x[length(x+1)/2])
25      }
26    }else if(quartile_type == 'one'){
27      position <- (length(x) + 1)/4
28      if (floor(position) == ceiling(position)){
29        return(x[position])
30      }else{
31        return((x[floor(position)] + x[ceiling(position)])/2)
32      }
33    }else if(quartile_type == 'three'){
34      position <- (3*(length(x) + 1))/4
35      if (floor(position) == ceiling(position)){
36        return(x[position])
```

```
37      }else{
38        return((x[floor(position)] + x[ceiling(position)])/2)
39      }
40    }
41 }
42
43 #Sample mean
44 cal_mean <- function(x){
45   result <- sum(x) / length(x)
46   return(result)
47 }
48
49 #Sample five number summary
50 cal_five_num_summary <- function(x){
51   min <- cal_minimum(x)
52   q1 <- cal_quartile(x, 'one')
53   q2 <- cal_quartile(x, 'two')
54   q3 <- cal_quartile(x, 'three')
55   max <- cal_maximum(x)
56   return(c(min, q1, q2, q3, max))
57   # return(list(min, q1, q2, q3, max))
58 }
59
60 #Sample Skewness
61 cal_coef_skew <- function(x){
62   s <- sqrt(cal_sample_variance(x))
63   mean <- cal_mean(x)
64   total <- 0
65   for(i in x){
66     total = total + (i - mean)^3
67   }
68
69   return(total/(s^3 * length(x)))
70
71 }
72
73 #Sample Kurtosis
74 cal_coef_kurt <- function(x){
75   s <- sqrt(cal_sample_variance(x))
76   mean <- cal_mean(x)
77   total <- 0
78   for(i in x){
79     total = total + (i - mean)^4
80   }
81   total = total/(s^4 * length(x))
82   return(total -3)
83 }
84
85 #Sample Variance
```

```
 86  cal_sample_variance <- function(x){
 87    mean <- cal_mean(x)
 88    n_minus_one <- length(x) - 1
 89    total <- 0
 90    for(i in x){
 91      total = total + (i - mean)^2
 92    }
 93    return(total/n_minus_one)
 94  }
 95
 96  #Sample Standard Deviation
 97  cal_sd <- function(x){
 98    sample_variance <- cal_sample_variance(x)
 99    return(sqrt(sample_variance))
100  }
101
102  #Sample Range
103  cal_sample_range <- function(x){
104    return(cal_maximum(x) - cal_minimum(x))
105  }
106
107  #Sample Decile Range
108  cal_sample_decile_range <- function(x){
109    x_0.9 <- 90/100 * (length(x) +1)
110    x_0.1 <- 10/100 * (length(x) +1)
111
112    if (floor(x_0.9) == ceiling(x_0.9)){
113      x_0.9 = x_0.9
114    }else{
115      x_0.9 = (x[floor(x_0.9)] + x[ceiling(x_0.9)])/2
116    }
117
118    if (floor(x_0.1) == ceiling(x_0.1)){
119      x_0.1 = x_0.1
120    }else{
121      x_0.1 = (x[floor(x_0.1)] + x[ceiling(x_0.1)])/2
122    }
123    return(x_0.9 - x_0.1)
124  }
125
126  #Sample 0.1 trimmed average
127  cal_sample_0.1_trimmed_average <- function(x){
128    g <- abs(0.1 * length(x))
129    n <- length(x)
130    total <- 0
131    for(i in seq(g+1, n-g, 1)){
132      total <- total + x[i]
133    }
134    return(total/(n - 2*g))
```

```
135  }
136
137  #Sample 0.1 trimmed variance
138  cal_sample_0.1_trimmed_variance <- function(x){
139     g <- abs(0.1 * length(x))
140     n <- length(x)
141     x_gt <- cal_sample_0.1_trimmed_average(x)
142     total <- 0
143     for(i in seq(g+1, n-g, 1)){
144        total <- total + (x[i] - x_gt)^2
145     }
146     return(total/(n - 2*g - 1))
147  }
148
149  library(xtable)
```

2. Separately for each population calculate sample size and these characteristics of maximum cranial breadth. Print them in a table with values rounded to 4 decimal places.

## Implementation in R

```
150  setwd('/home/jane/Documents/SSME (Service Development Management)/Term2/
        MV013-Statistics for Computer Science/Assignment1')
151  Howell <- read.csv('Howell.csv')
152  Howell.Male <- Howell[Howell$Sex == 'M', ]
153  Howell.Male.Australi <- Howell.Male[(Howell.Male$Population == 'AUSTRALI
        '), ]
154  Howell.Male.Peru <- Howell.Male[(Howell.Male$Population == 'PERU'), ]
155
156  #Working on population of 'AUSTRALI'
157  Australi.sample_size <- length(Howell.Male.Australi$XCB)
158  Australi.sample_mean <- cal_mean(Howell.Male.Australi$XCB)
159  Australi.sample_five_num_summary <- cal_five_num_summary(Howell.Male.
        Australi$XCB)
160  Australi.sample_skewness <- cal_coef_skew(Howell.Male.Australi$XCB)
161  Australi.sample_kurtosis <- cal_coef_kurt(Howell.Male.Australi$XCB)
162  Australi.sample_variance <- cal_sample_variance(Howell.Male.Australi$XCB
        )
163  Australi.sample_sd <- cal_sd(Howell.Male.Australi$XCB)
164  Australi.sample_range <- cal_sample_range(Howell.Male.Australi$XCB)
165  Australi.sample_decile_range <- cal_sample_decile_range(Howell.Male.
        Australi$XCB)
166  Australi.sample_0.1_trimmed_average <- cal_sample_0.1_trimmed_average(
        Howell.Male.Australi$XCB)
167  Australi.sample_0.1_trimmed_variance <- cal_sample_0.1_trimmed_variance(
        Howell.Male.Australi$XCB)
168
169  #Working on population of 'PERU'
170  Peru.sample_size <- length(Howell.Male.Peru$XCB)
171  Peru.sample_mean <- cal_mean(Howell.Male.Peru$XCB)
```

```
172  Peru.sample_five_num_summary <- cal_five_num_summary(Howell.Male.Peru$
        XCB)
173  Peru.sample_skewness <- cal_coef_skew(Howell.Male.Peru$XCB)
174  Peru.sample_kurtosis <- cal_coef_kurt(Howell.Male.Peru$XCB)
175  Peru.sample_variance <- cal_sample_variance(Howell.Male.Peru$XCB)
176  Peru.sample_sd <- cal_sd(Howell.Male.Peru$XCB)
177  Peru.sample_range <- cal_sample_range(Howell.Male.Peru$XCB)
178  Peru.sample_decile_range <- cal_sample_decile_range(Howell.Male.Peru$XCB
        )
179  Peru.sample_0.1_trimmed_average <- cal_sample_0.1_trimmed_average(Howell
        .Male.Peru$XCB)
180  Peru.sample_0.1_trimmed_variance <- cal_sample_0.1_trimmed_variance(
        Howell.Male.Peru$XCB)
181
182  Peru <- data.frame(Peru=c(Peru.sample_size, Peru.sample_mean, Peru.
        sample_five_num_summary, Peru.sample_skewness,
183                          Peru.sample_kurtosis, Peru.sample_variance, Peru
                                .sample_sd, Peru.sample_range,
184                          Peru.sample_decile_range, Peru.sample_0.1_
                                trimmed_average, Peru.sample_0.1_trimmed_
                                variance),
185                   row.names = c('Size', 'Mean', 'Minimum', 'Q1', 'Q2', '
                          Q3', 'Maximum', 'Skewness', 'Kurtosis', 'Variance',
186                          'Standard Deviation', 'Range', 'Decile Range', '
                                Sampe 0.1 trimmed average', 'Sample 0.1
                                trimmed variance'))
187
188  Australi <- data.frame(Australi=c(Australi.sample_size, Australi.sample_
        mean, Australi.sample_five_num_summary, Australi.sample_skewness,
189                          Australi.sample_kurtosis, Australi.sample_
                                variance, Australi.sample_sd, Australi.sample
                                _range,
190                          Australi.sample_decile_range, Australi.sample_
                                0.1_trimmed_average, Australi.sample_0.1_
                                trimmed_variance),
191                       row.names = c('Size', 'Mean', 'Minimum', 'Q1', '
                              Q2', 'Q3', 'Maximum', 'Skewness', 'Kurtosis',
                              'Variance',
192                          'Standard Deviation', 'Range', 'Decile Range', '
                                Sampe 0.1 trimmed average', 'Sample 0.1
                                trimmed variance'))
193
194  library(xtable)
```

3. Create boxplots of maximum cranial breadth of each population side by side in one figure. Set the width of boxes to be proportional to sample sizes, add notches and arithmetic averages for both groups.

**Implementation in R**

```
195  max.len = max(length(Howell.Male.Australi$XCB), length(Howell.Male.Peru$
         XCB))
196  x=c(Howell.Male.Australi$XCB, rep(NA, max.len - length(Howell.Male.
         Australi$XCB)))
197  y=c(Howell.Male.Peru$XCB, rep(NA, max.len - length(Howell.Male.Peru$XCB)
         ))
198  Populations <- c('Australi', 'Peru')
199  values <- c(x, y)
200  df <- data.frame(Populations, values)
201
202  library(ggplot2)
203  boxplot <- ggplot(df, aes(x=Populations, y=values, fill=Populations)) +
204            geom_boxplot(notch = TRUE, width=(55/110)) +
205            scale_fill_manual(values=c("lightblue", "lightgreen")) +
206            ylab('Maximum cranial breadth') +
207            xlab('Populations') +
208            annotate("text", x=1, y=Australi.sample_mean, label= round(
                 Australi.sample_mean, 4)) +
209            annotate("text", x=2, y=Peru.sample_mean, label= round(Peru.
                 sample_mean, 4)) +
210            ggtitle('Boxplot of maximum cranial breadth') +
211            theme(plot.title = element_text(hjust = 0.5))
```

4. Separately for each population create histogram of maximum cranial breadth. Make sure that the histograms can be easily compared (without using back-to-back histogram).

## Implementation in R

```
212  Australi.hist <- ggplot(data=data.frame(Australi=c(Howell.Male.Australi$
         XCB)), aes(x=Howell.Male.Australi$XCB)) +
213                geom_histogram(color="darkblue", fill="lightblue",
                     binwidth = 1) +
214                xlab('Maximum cranial breadth of Australi') +
215                ylab('Frequecy') +
216                scale_x_continuous(breaks = seq(124, 144, 2), lim = c
                     (122, 146)) +
217                geom_text(stat = 'count', aes(label =..count.., vjust =
                     -0.5))
218
219  Peru.hist <- ggplot(data=data.frame(Peru=c(Howell.Male.Peru$XCB)), aes(x
         =Howell.Male.Peru$XCB)) +
220            geom_histogram(color="darkgreen", fill="lightgreen",
                 binwidth = 1) +
221            xlab('Maximum cranial breadth of Peru') +
222            ylab('Frequency') +
223            scale_x_continuous(breaks = seq(129, 149, 2), lim = c(128,
                 150)) +
224            geom_text(stat = 'count', aes(label =..count.., vjust =
                 -0.5))
```

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

## Implementation in R

```
225  Australi.qq <- ggplot(data=data.frame(Australi=c(Howell.Male.Australi$
        XCB))) +
226                  aes(sample= Australi) +
227                  geom_qq(distribution = qnorm) +
228                  geom_qq_line(line.p = c(0.25, 0.75), col = "lightblue")
                        +
229                  ylab('Maximum cranial breadth for Australi') +
230                  ggtitle('QQ-Plot of maximum cranial breadth of Australi'
                        ) +
231                  theme(plot.title = element_text(hjust = 0.5))
232
233  Peru.qq <- ggplot(data=data.frame(Peru=c(Howell.Male.Peru$XCB))) +
234                aes(sample= Peru) +
235                geom_qq(distribution = qnorm) +
236                geom_qq_line(line.p = c(0.25, 0.75), col = "lightgreen") +
237                ylab('Maximum cranial breadth for Peru') +
238                ggtitle('QQ-Plot of maximum cranial breadth of Peru') +
239                theme(plot.title = element_text(hjust = 0.5))
```
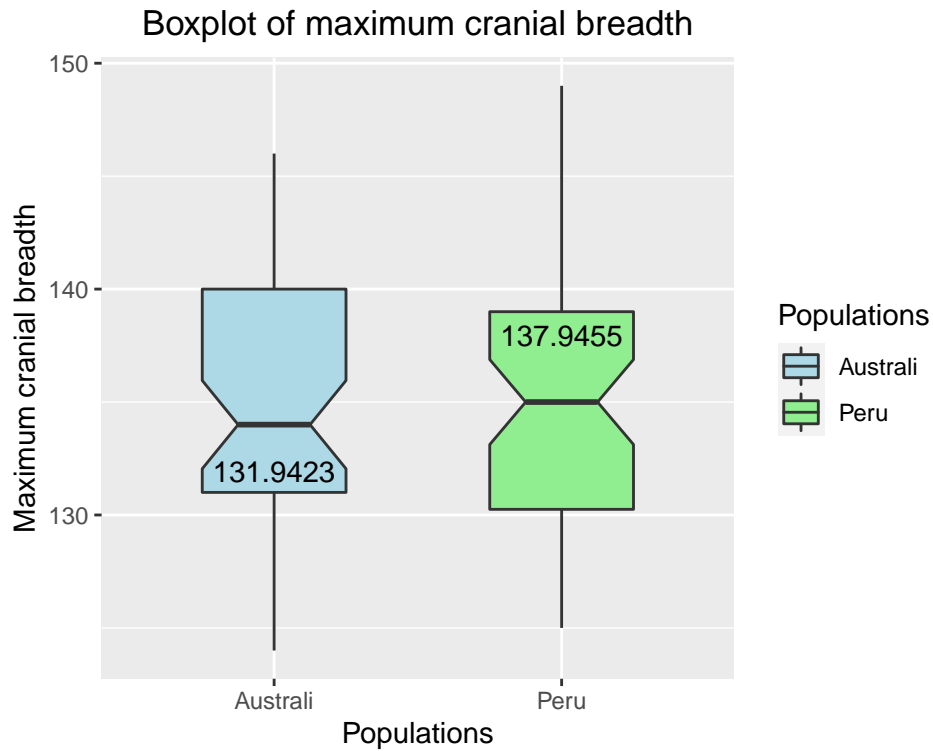
6. Interpret your results and graphics. Text. Results in table or graphic form. Commentaries and interpretation of the results.
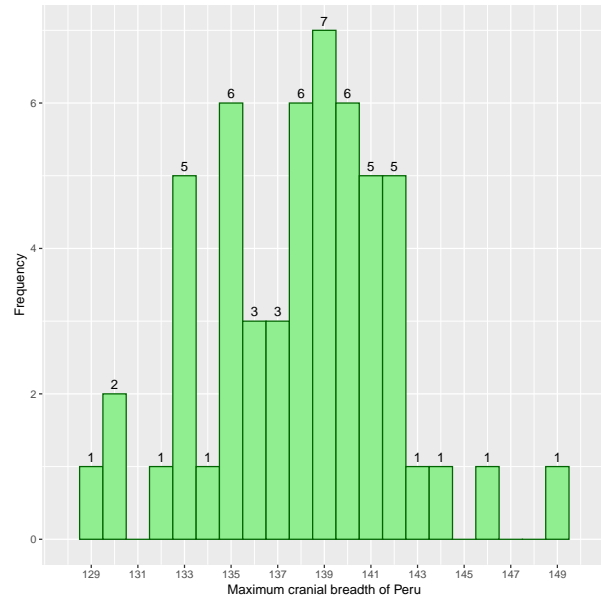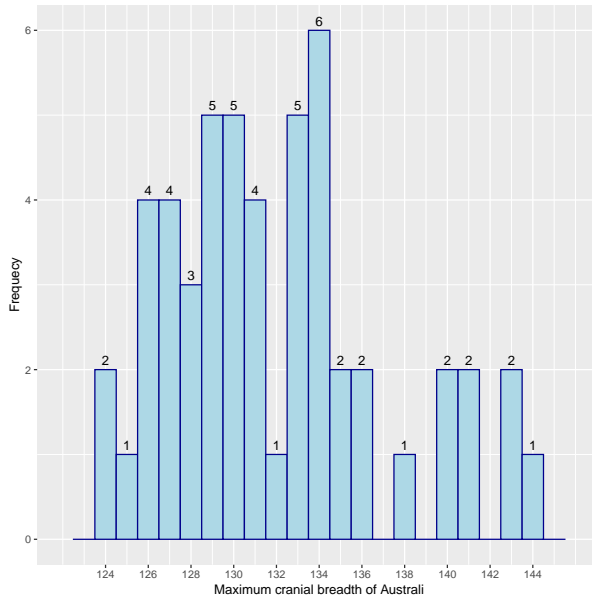
### Characteristics of maximum cranial breadth

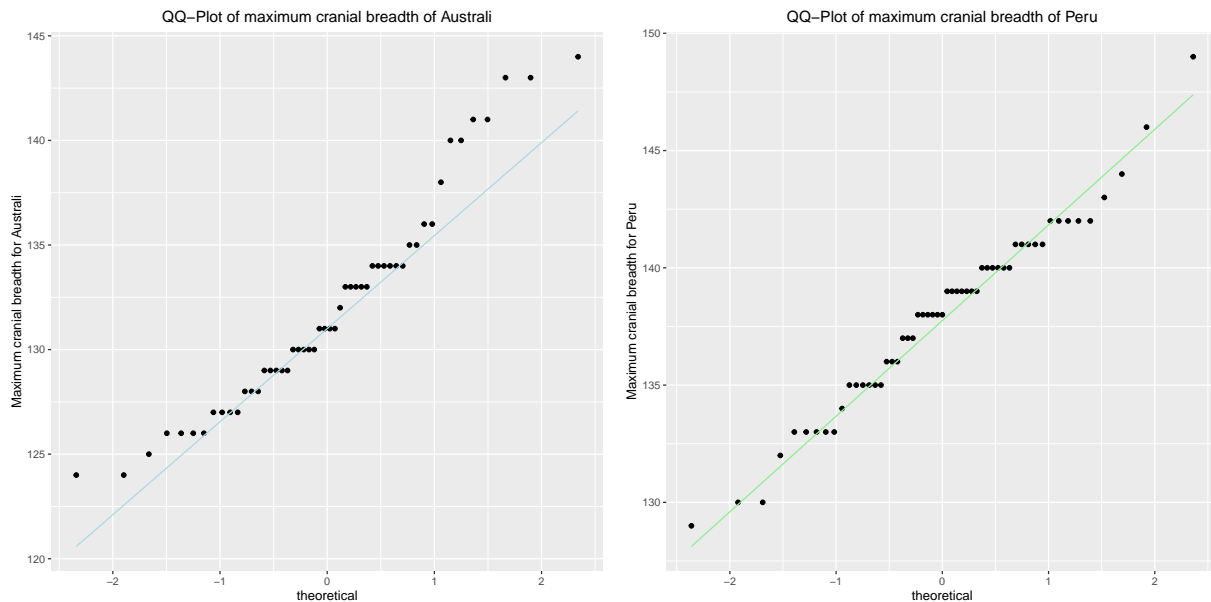| | Australi | | Peru |
|---|---|---|---|
| Size | 52.0000 | Size | 55.0000 |
| Mean | 131.9423 | Mean | 137.9455 |
| Minimum | 124.0000 | Minimum | 129.0000 |
| Q1 | 128.0000 | Q1 | 135.0000 |
| Q2 | 131.0000 | Q2 | 138.0000 |
| Q3 | 134.0000 | Q3 | 141.0000 |
| Maximum | 144.0000 | Maximum | 149.0000 |
| Skewness | 0.6436 | Skewness | -0.0037 |
| Kurtosis | -0.3390 | Kurtosis | 0.0616 |
| Variance | 26.0554 | Variance | 15.8673 |
| Standard Deviation | 5.1045 | Standard Deviation | 3.9834 |
| Range | 20.0000 | Range | 20.0000 |
| Decile Range | 15.0000 | Decile Range | -4.0000 |
| Sampe 0.1 trimmed average | 129.4712 | Sampe 0.1 trimmed average | 138.1364 |
| Sample 0.1 trimmed variance | 26.3226 | Sample 0.1 trimmed variance | 15.1438 |

## Boxplot of maximum cranial breadth



### Histograms represent maximum cranial breadth of each population

QQ-plot represent maximum cranial breadth of each population



# Exercise 2

1. Calculate the number of men and women in Spain for each year and print them in a table together with the total population.

## Implementation in R

```
240  setwd('/home/jane/Documents/SSME (Service Development Management)/Term2/
        MV013-Statistics for Computer Science/Assignment1')
241  Spanish.province <- read.csv('area_spanish_provinces.csv')
242  Spanish.population <- read.csv('population-spain-1998-2018.csv', sep=";"
        )
243
244  population.2018 <- sum(Spanish.population$males.2018, Spanish.population
        $females.2018)
245  population.2013 <- sum(Spanish.population$males.2013, Spanish.population
        $females.2013)
246  population.2008 <- sum(Spanish.population$males.2008, Spanish.population
        $females.2008)
247  population.2003 <- sum(Spanish.population$males.2003, Spanish.population
        $females.2003)
248  population.1998 <- sum(Spanish.population$males.1998, Spanish.population
        $females.1998)
249
250  total.population <- sum(colSums(Spanish.population[ ,-1]))
251
252  df <- data.frame(Male=c(sum(Spanish.population$males.1998), sum(Spanish.
        population$males.2003),
```

```
253                          sum(Spanish.population$males.2008), sum(Spanish.
                                 population$males.2013),
254                          sum(Spanish.population$males.2018), 0),
255              Female=c(sum(Spanish.population$females.1998), sum(
                         Spanish.population$females.2003),
256                          sum(Spanish.population$females.2008), sum(
                                 Spanish.population$females.2013),
257                          sum(Spanish.population$females.2018), 0),
258              Total=c(population.1998, population.2003, population
                         .2008, population.2013, population.2018, total.
                         population))
259 row.names(df) <- c('1998', '2003', '2008', '2013', '2018', 'Total')
260 df[6,1] <- sum(df$Male)
261 df[6,2] <- sum(df$Female)
```

2. Display barplot plot of total population of Spain in each of the years, with each bar divided
between men and women.

## Implementation in R

```
262 library(ggplot2)
263 year <- c(1998, 1998, 2003, 2003, 2008, 2008, 2013, 2013, 2018, 2018)
264 population <- c(sum(Spanish.population$males.1998),
265                 sum(Spanish.population$females.1998),
266                 sum(Spanish.population$males.2003),
267                 sum(Spanish.population$females.2003),
268                 sum(Spanish.population$males.2008),
269                 sum(Spanish.population$females.2008),
270                 sum(Spanish.population$males.2013),
271                 sum(Spanish.population$females.2013),
272                 sum(Spanish.population$males.2018),
273                 sum(Spanish.population$females.2018))
274 gender <- c('Male', 'Female', 'Male', 'Female', 'Male', 'Female', 'Male'
       , 'Female', 'Male', 'Female' )
275 df.population <- data.frame(year, population, gender)
276 population.barplot<-ggplot(data=df.population, aes(x=year, y=population,
       fill=gender )) +
277   geom_col() +
278   scale_fill_manual(values=c("lightpink", "lightblue")) +
279   scale_x_continuous(breaks = seq(1998, 2018, 5)) +
280   # scale_y_continuous(breaks = sprintf("%.0fk", population/1000)) +
281   geom_text(aes(label = sprintf("%.0fk", population/1000)), position =
       position_stack(0.5)) +
282   theme(axis.text.y=element_text(angle = 45))
```

3. Display barplot of relative proportions of men and women within each province in 2018

## Implementation in R

```
283 proportion.province <-c(as.vector(Spanish.population$province), as.
       vector(Spanish.population$province))
```

```
284  proportion.gender <- c(rep('Male', 52), rep('Female', 52))
285  proportion.province.gender <- c(proportion.male, proportion.female)
```

```
Error in eval(expr, envir, enclos): object 'proportion.male' not found     286
```

```
287  df.proportion <- data.frame(proportion.province, proportion.province.
         gender, proportion.gender)
```

```
Error in data.frame(proportion.province, proportion.province.gender,       288
    proportion.gender): object 'proportion.province.gender' not found
```

```
289  proportion.barplot <- ggplot(data=dffff, aes(x=proportion.province, y=
         proportion.province.gender, fill=proportion.gender)) +
290          geom_col() +
291          theme(axis.text.x=element_text(angle = 90, size=6), plot.title
                 = element_text(vjust=0.5, hjust=1)) +
292          scale_fill_manual(values=c("lightpink", "lightblue")) +
293          scale_x_discrete(name= "Provinces") +
294          scale_y_continuous(name= "Population proportion") +
295          ggtitle('Barplot of population proportion')
```

```
Error in ggplot(data = dffff, aes(x = proportion.province, y =             296
    proportion.province.gender, : object 'dffff' not found
```

```
297  proportion.barplot
```

```
Error in eval(expr, envir, enclos): object 'proportion.barplot' not        298
    found
```

5. Interpret your results and graphics.

| | Male | Female | Total |
|---|---|---|---|
| 1998 | 19488465 | 20364186 | 39852651 |
| 2003 | 21034326 | 21682738 | 42717064 |
| 2008 | 22847737 | 23310085 | 46157822 |
| 2013 | 23196386 | 23933397 | 47129783 |
| 2018 | 22896602 | 23826378 | 46722980 |
| Total | 109463516 | 113116784 | 222580300 |

Table 1: Population of Spanish