

Assignment 1

NAME: Nicholas Tolley

DUE DATE: February 14th, 6pm

Please post your completed homework (including the R code).

Problem I (100 pts)

This exercise focuses on the data sets `country.rda`, which can be found in Canvas in the Files item in the folder Data (file name: `country.rda`). It contains data called `gapminder` from the Gapminder non-profit organization. The data are 10,545 units and 9 variables. The variables are:

`country`

`continent`

`year`

`infant_mortality`: infant mortality rate (to interpret, if you divide by 10 you will obtain the percentage)

`life_expectancy`

`fertility`

`population`

`gdp`

`continent`

`region`

Read and load the data in R.

- 1) Perform a boxplot with `ggplot` in R where you can see in the x axis each continent and y axis the population just for year 1960, hint: use `population/10^6` and use the log transformation. Be careful on writing the axes labels. What continent has the largest median population size? Check in R. Explain what you have found with theoretical details (outliers, symmetry of the distribution, median and interquartile range).
- 2) Perform a boxplot with `ggplot` where you can see in the x axis each continent and y axis the population just for year 2010, hint: use `population/10^6` and use the log transformation. What continent has the largest median population size in this year? Check in R. Explain what you have found with theoretical details (outliers, symmetry of the distribution, median and interquartile range).
- 3) What is the median population size for Africa to the nearest million? What proportion of countries in Europe have populations below 14 million? Perform in R.
- 4) If we use a log transformation, which continent has the largest interquartile range? Perform and plot in R. Explain what you have found, give also the theoretical implications.

Compare (in `ggplot`) the countries in term of population using boxplots with two different colours for year 1960 and 2010. Try to make a plot exactly as Figure 1 (just the colour can change).

- 6) Compare (in `ggplot`) the countries in term of infant mortality using boxplots with two different colours for year 1960 and 2010. Try to make a plot exactly as Figure 2 (just the colour can change).

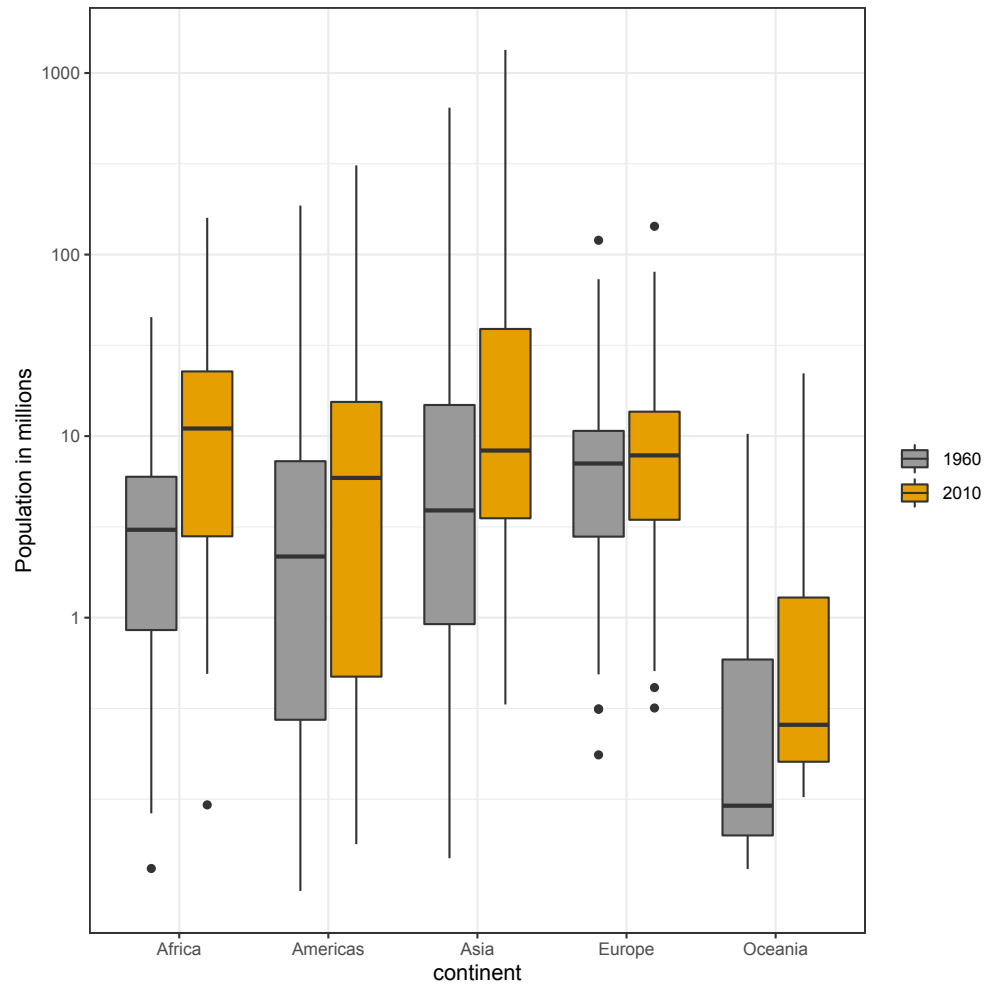


Figure 1: Population per country

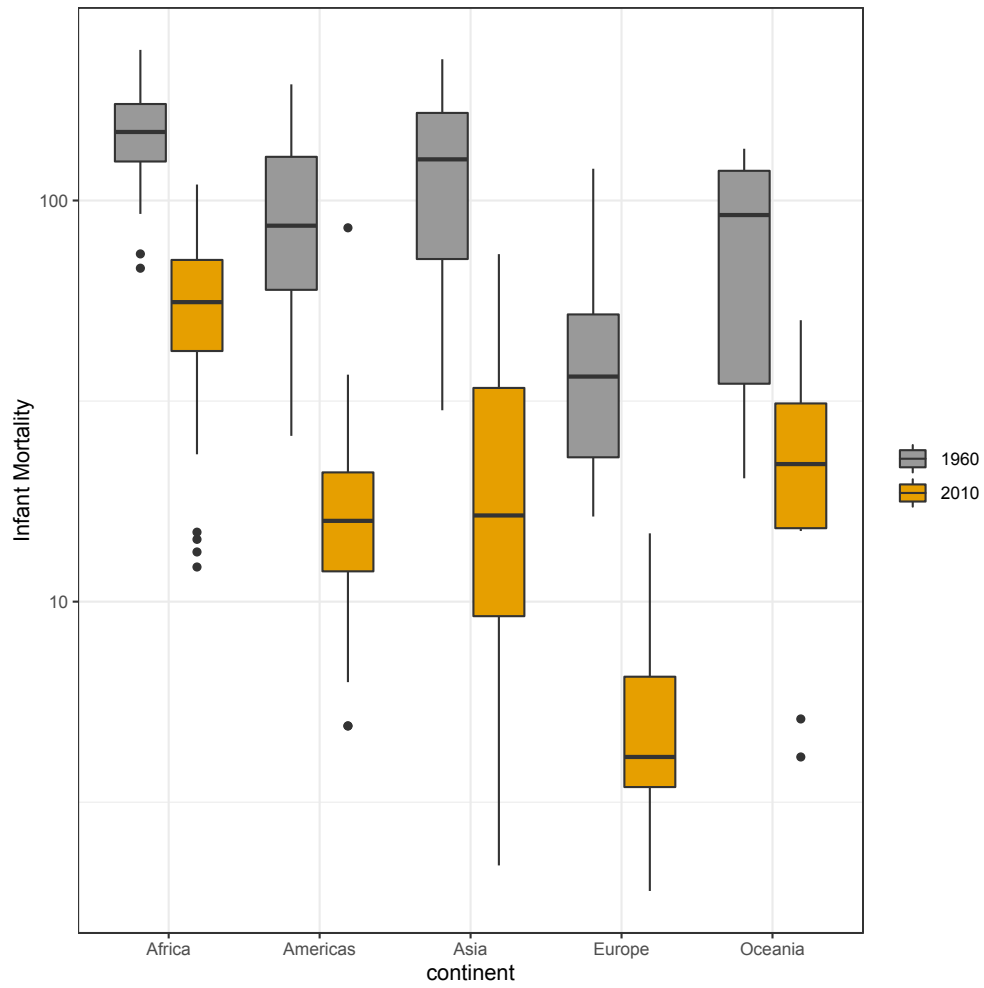


Figure 2: Infant Mortality per country

- 7) What proportion of the infant mortality is between 100 and 150? And what is its approximation to a normal distribution of this proportion? What proportion of the infant mortality is greater than 150? And what is its approximation to a normal distribution of the proportions? Perform in R. Explain what you have found, give also the theoretical implications and what theorem or statistical assumptions you are using.
- 8) In ggplot perform both a histogram and a curve for infant mortality, use two different colours for year 1960 and 2010. Does the infertility follow a normal distribution? Are the two distribution symmetric or skewed? Be careful to choose the right x axis for the histogram.
- 9) Perform a q-q plot for the infant mortality to check if this variables is normal or not. Write the label of the axes in the appropriate way. Explain what you have found including the theoretical implications.
- 10) Compare the regions in term of infant mortality for the two years 1960 and 2010 using the barplot. Try to make a plot exactly as Figure 3 (just the colour can change). Create a binary variable, with value 1 if life-expectancy is greater than 65 (group A), 0 otherwise (group B). Then perform the appropriate two sample test to know if the mean of the infant mortality is the same in group A and B (remember to check the variance of the two groups). Explain your results.

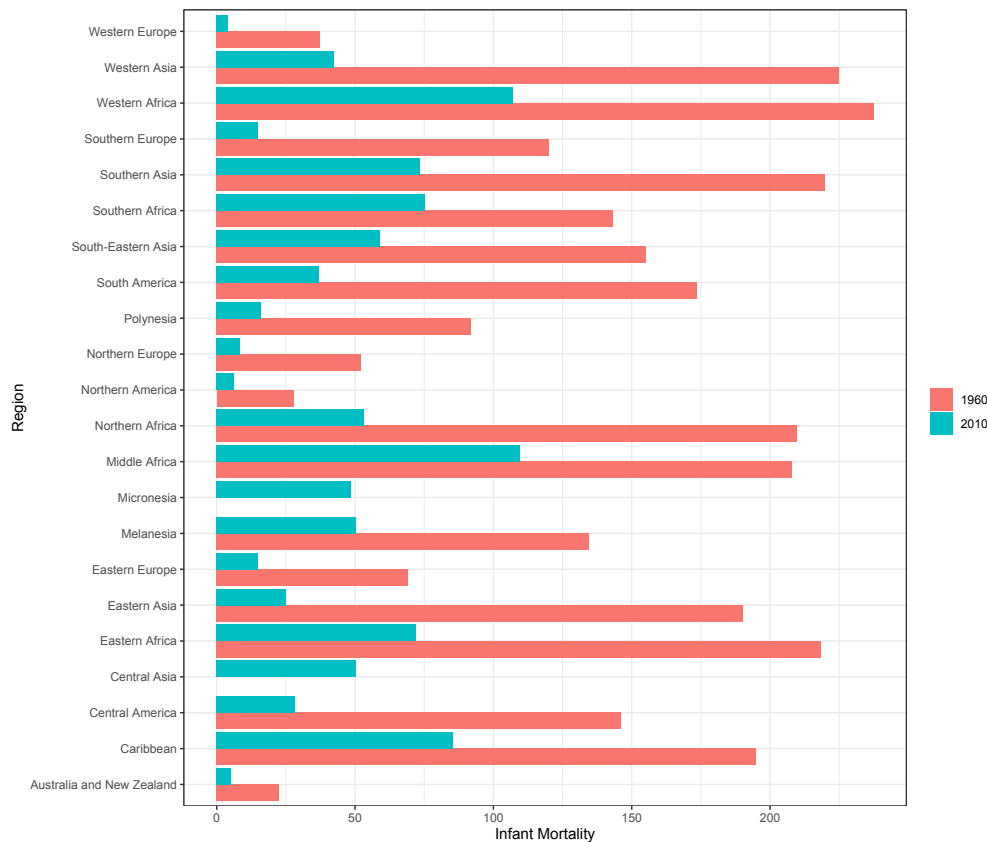


Figure 3: Infant Mortality per region