

# Gapminder Example Report

*Melinda Higgins*

*January 16, 2017*

## Gapminder Example

The following report is based on the R code examples at the **Gapminder** Github repository located at <https://github.com/jennybc/gapminder>.

## Gapminder dataset - summary statistics

The `gapminder` dataset which is built into the `gapminder` package, has 6 variables and 1704 observations. We can list the variables in the dataset, using the `names()` function.

### Variables (columns) in gapminder dataset

```
names(gapminder)

## [1] "country" "continent" "year" "lifeExp" "pop" "gdpPercap"
```

### Structure of the gapminder dataset

Another way to see the “structure” of the dataset is to run the `str()` function.

```
str(gapminder)

## Classes 'tbl_df', 'tbl' and 'data.frame': 1704 obs. of 6 variables:
## $ country : Factor w/ 142 levels "Afghanistan",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ continent: Factor w/ 5 levels "Africa","Americas",...: 3 3 3 3 3 3 3 3 3 3 ...
## $ year : int 1952 1957 1962 1967 1972 1977 1982 1987 1992 1997 ...
## $ lifeExp : num 28.8 30.3 32 34 36.1 ...
## $ pop : int 8425333 9240934 10267083 11537966 13079460 14880372 12881816 13867957 16317921 22...
## $ gdpPercap: num 779 821 853 836 740 ...
```

You’ll notice that the 1st 2 columns/variables “country” and “continent” are both “Factor” type variables. Variables 3 and 5 “year” and “pop” are both “int” Integer type variables. Variables 4 and 6 “lifeExp” and “gdpPercap” are both “num” Numeric type variables.

### Summary Statistics of the gapminder dataset

The built in function `summary()` in base R does a good simple summary statistics for all variables in the dataset provided. Since this dataset only has 6 variable, we can simply call `summary(gapminder)` which will give us the summary statistics for all 6 variables.

```
summary(gapminder)

## country continent year lifeExp
## Afghanistan: 12 Africa :624 Min. :1952 Min. :23.60
## Albania : 12 Americas:300 1st Qu.:1966 1st Qu.:48.20
```

```
## Algeria      : 12   Asia      :396   Median :1980   Median :60.71
## Angola       : 12   Europe   :360   Mean    :1980   Mean    :59.47
## Argentina    : 12   Oceania : 24   3rd Qu.:1993   3rd Qu.:70.85
## Australia    : 12                   Max.    :2007   Max.    :82.60
## (Other)      :1632
##      pop      gdpPercap
## Min.      :6.001e+04   Min.      : 241.2
## 1st Qu.:2.794e+06   1st Qu.: 1202.1
## Median :7.024e+06   Median : 3531.8
## Mean    :2.960e+07   Mean    : 7215.3
## 3rd Qu.:1.959e+07   3rd Qu.: 9325.5
## Max.    :1.319e+09   Max.    :113523.1
##
```

### Run specific statistic for a given variable

Suppose we only wanted to get the *mean* life expectancy. To do this we can use the built-in function `mean()`. To select only the “lifeExp” variable, we can either refer to it by which column it is in the dataset using `gapminder[,4]` which says to select all rows by leaving the 1st element between the `[]` before the comma blank and putting a 4 after the comma which specifies the 4th column. Another way to select a column is to use the name of that column which is “lifeExp” and use the dollar sign `$` selector to get `gapminder$lifeExp`.

```
mean(gapminder$lifeExp)
```

```
## [1] 59.47444
```

### In line code

We can use the same command above, but call it “inline” instead of as a separate code chunk which sets the output apart from the text in a separate section of the report. If you simply want the computation executed and the result inserted into the body of text you are writing you use “inline” code. To do this you use the following syntax `r mean(gapminder$lifeExp)` between the backtick marks ‘```’. Using this syntax, we can write the following sentence.

The mean life expectancy is 59.4744394 years.

We can clean this up further by wrapping this command within the `round()` function and specifying the number of digits we want reported for this numeric result. This time, use the following command inline `r round(mean(gapminder$lifeExp), digits=2)` and we’ll rewrite the sentence below.

The mean life expectancy is 59.47 years.

### Homework 01 Exercise - Task 1

Modify the sentence above to also provide the standard deviation, median and sample size for life expectancy, set `digits=2`.

Hint: Read help pages for the functions `sd()`, `median()`, and `length()`.

### Look at a statistic by continent

Using the 1st code example at the `gapminder` Github repository at <https://github.com/jennybc/gapminder>, use the `aggregate()` command to see the median life expectancy by continent. You’ll notice that the 1st variable listed is the “lifeExp” variable we want run “by” “continent”. The “by” is indicated using the *tilde*

symbol ~. The 2nd variable listed is “continent” - this 2nd variable is usually a “factor-type” variable or group variable. *Hint: Try running lifeExp by year to get median lifeExp for each year.*

```
aggregate(lifeExp ~ continent, gapminder, median)
```

```
##   continent lifeExp
## 1   Africa 47.7920
## 2 Americas 67.0480
## 3   Asia 61.7915
## 4  Europe 72.2410
## 5  Oceania 73.6650
```

## Homework 01 Exercise - Task 2

Modify the r code chunk above to also provide the mean and standard deviation for life expectancy by continent.