

# VU BI2 - Exercise 1: Olympics Data

*Asil Cetin / 01100130*

*6/11/2017*

## Data Description

In this exercise we are investigating a data set from London 2012 Olympics. The data set gives the names of the 203 participating countries as well as the number of gold, silver and bronze medals won by country, the total number of medals won by country, the Borda points by country, income per capita (in \$1.000), population size (in 1.000.000), gross domestic product (GDP= income per capita multiplied by population size) and the polynomial variables of income per capita squared, population size squared, income per capita cubed, population size cubed, gross domestic product squared, gross domestic product cubed, natural log of income per capita, natural log of population size, and natural log of GDP.

## Analysis Questions

We are mainly interested in the correlation between the overall success in London 2012 Olympics - which is represented in the parameter “BordaPoints”, since it ranks the countries weighted on the value of different medals - and the population and income levels of a given country.

Thus our first analysis question can be stated as:

**Do parameters of population size and income have any effect on success in the London 2012 Olympics?**

After investigating the possible correlations between olympic success and parameters of population size and income, we would want to know to what degree these parameters have an effect. Thus our second question would be:

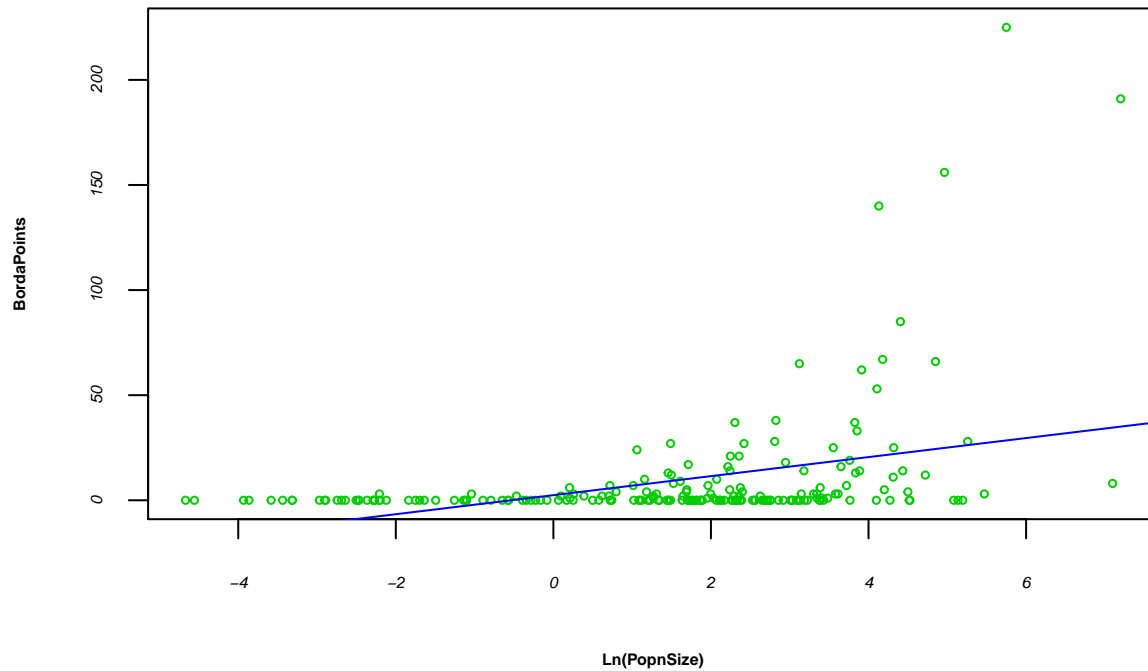
**If population size and income has an effect on olympic success, what are the factoring weights of these parameters?**

## Population and Success

First parameter we want to analyse is the population. Here we will use the natural logarithm of population size (in 1.000.000) since it's preferred to naturalize the enormous population size differences between some countries.

The plot of all participating countries in the 2012 London Olympics looks like as follows:

*Distribution of Borda Points by Population (ln)*

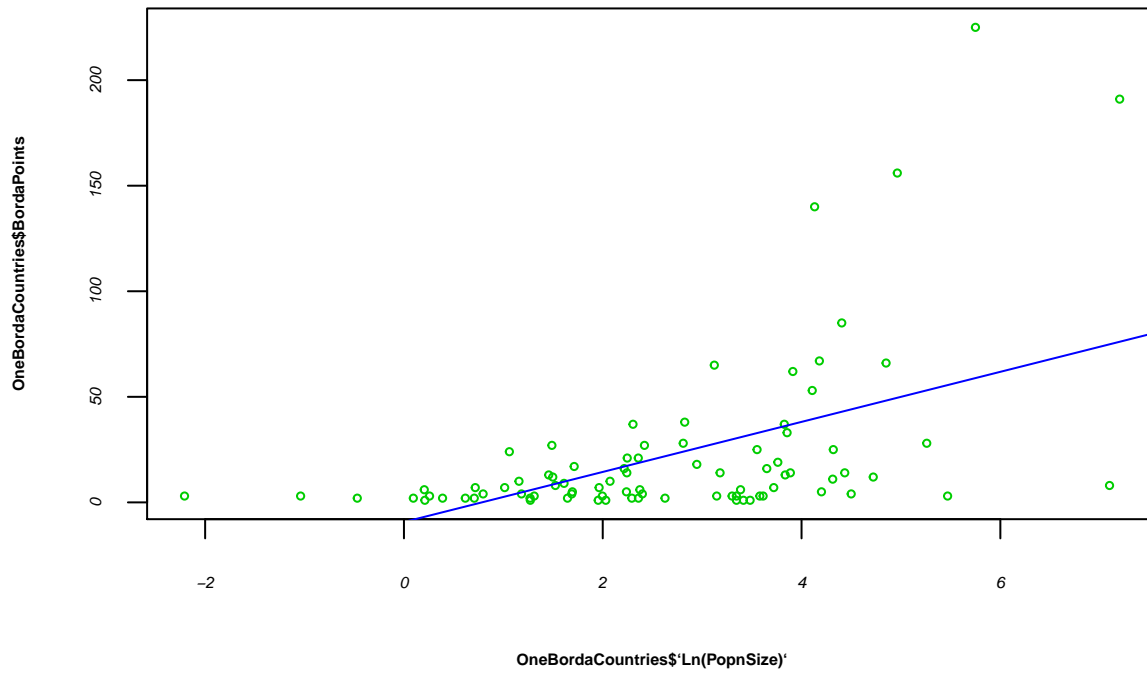


Population and success at the olympics don't show a strict correlation at the first glance.

However one should not forget that we observe many countries - with big and small populations - having no success (0 Borda Points), which may hinder us coming to effective conclusions.

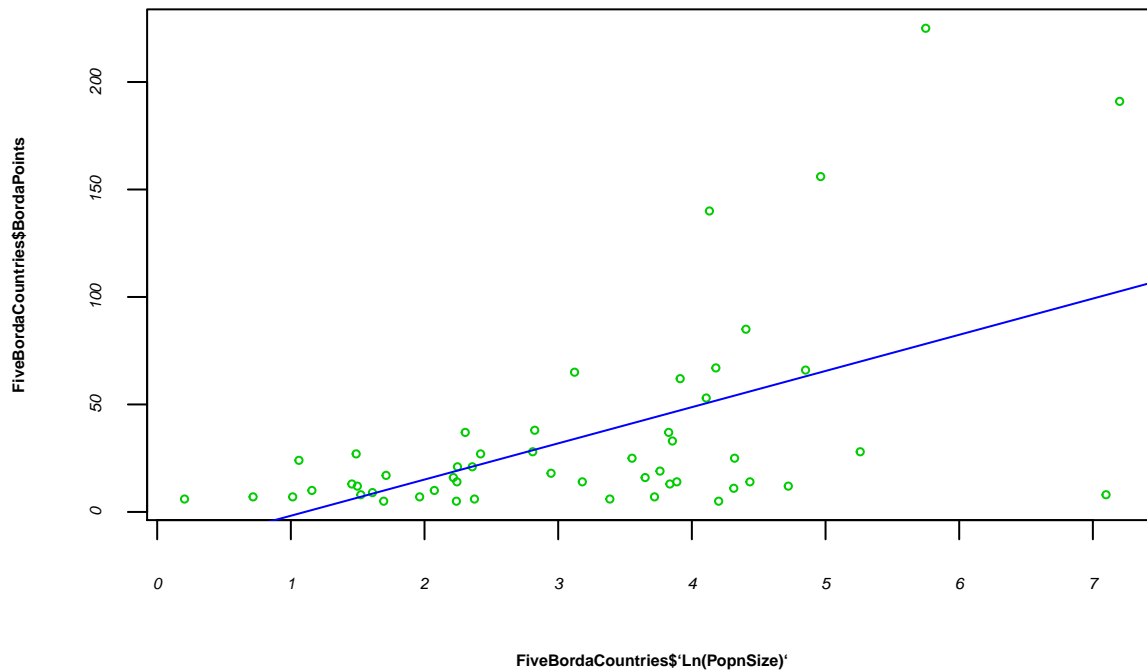
How would the same graphic look, if we only take countries with at least 1 Borda Point into account? This may be a better comparison since those are the countries which showed at least a minimal level of competitive participation in the events. The following plot depicts those countries with at least 1 Borda Point:

*Distribution of Borda Points by Population (ln) of Countries at least 1 Borda Point*



If we would like to further increase our minimal requirement of Borda Points from 1 upto 5, it might be argued that this would lead to a more precise analysis, since getting one Borda Point is done by only one bronze medal, which would be too small of a difference between 0 and 1 Borda Point countries. The new plot would look like the following:

*Distribution of Borda Points by Population (ln) of Countries at least 5 Borda Point*

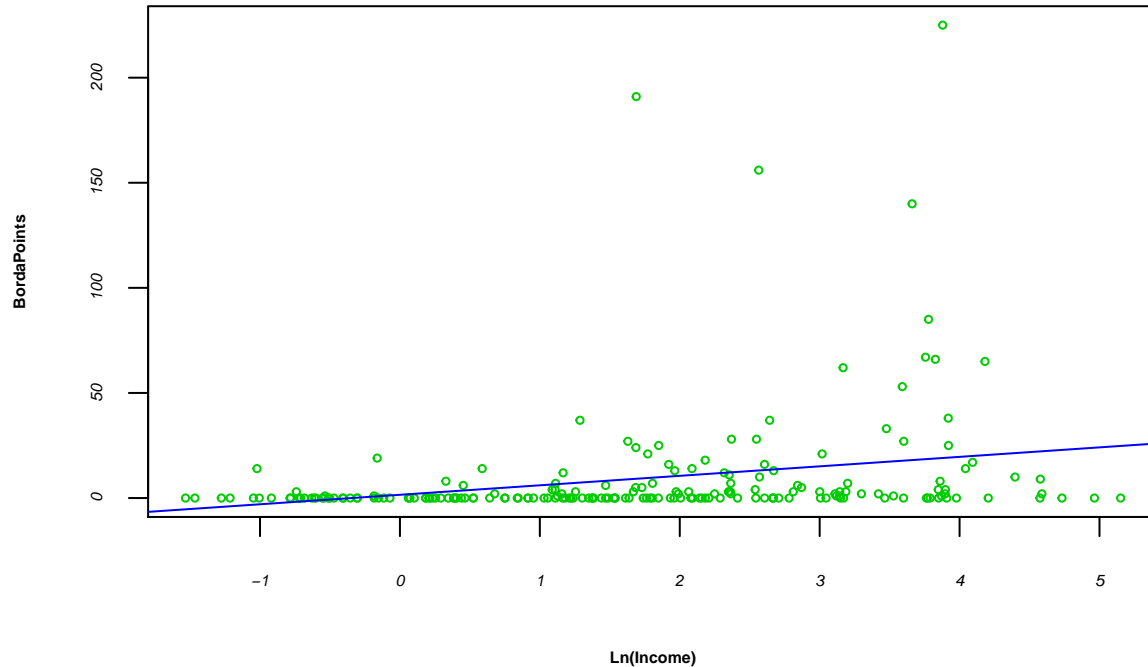


## Income and Success

Second parameter we want to analyse is the level of income. Here we will use the natural logarithm of income per capita (in \$1.000) since it's preferred to naturalize the enormous income size differences between some countries. Moreover income per capita parameter is preferred over GDP of a country, because we're interested in individuals' financial opportunities rather than a country's total production.

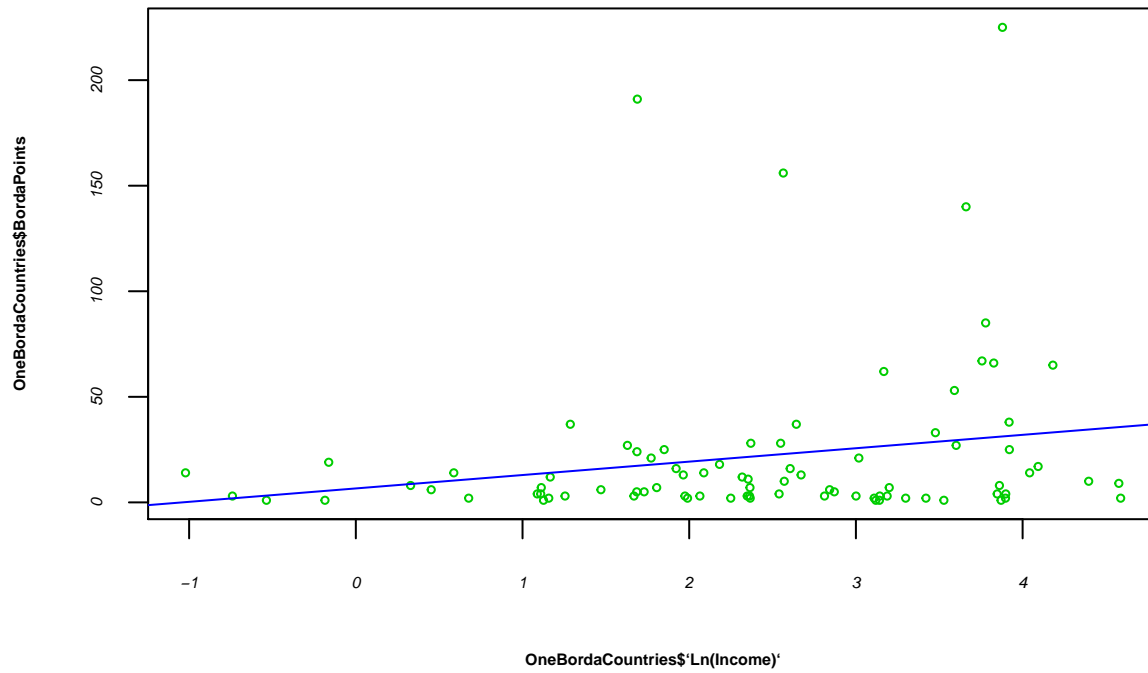
The plot of all participating countries in the 2012 London Olympics looks like as follows:

*Distribution of Borda Points by Income (ln)*

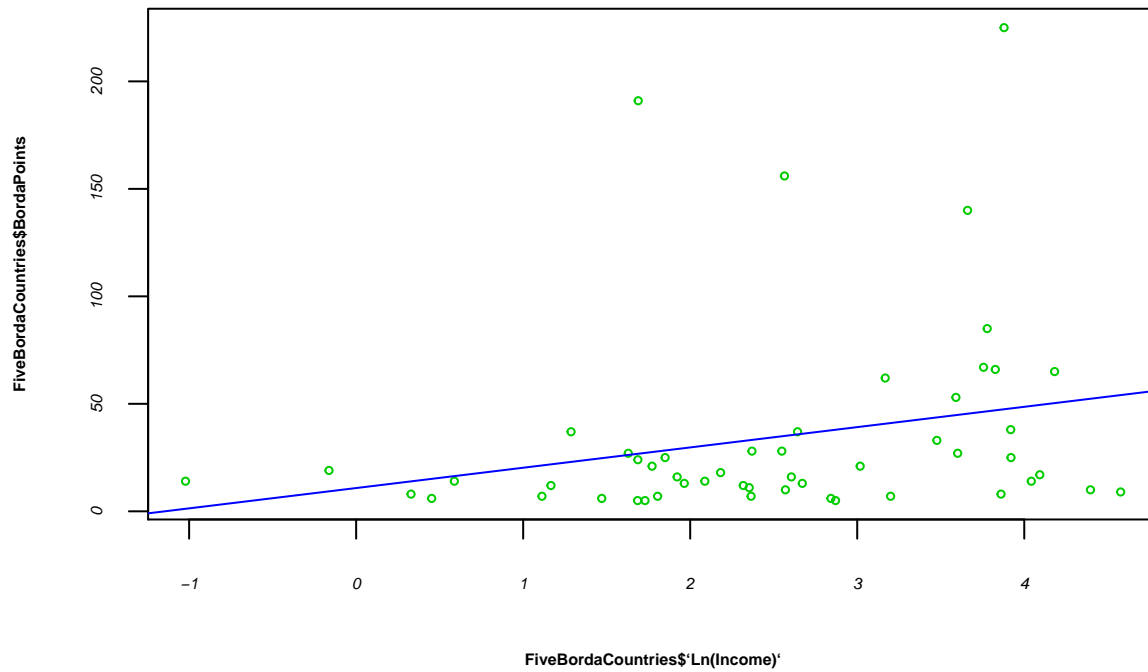


Similar to the first plot of population size vs. success, this plot doesn't represent a direct correlation due to a high number of countries with 0 Borda Points. Let's increase our Borda Point requirement to at least one and then to five points:

*Distribution of Borda Points by Income (ln) of Countries at least 1 Borda Point*



*Distribution of Borda Points by Income (ln) of Countries at least 5 Borda Point*



Our first descriptive analysis on both income and population have shown a certain degree of correlation with success, especially after restricting our data to countries which have achieved at least minimal level of medals in the competitions.

## Top 50 and Bottom 50 Countries and Averages

In this section we'd like to compare the summaries of (1) all participating countries, (2) top-50 countries by Borda Points and (3) bottom-50 countries by Borda Points.

First parameter to be analysed is the population size (in in 1.000.000):

Summary of all countries' population size:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.0094	1.2306	6.4980	34.0245	22.6124	1340.0865

Summary of top-50 countries' population size:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	1.227	8.256	23.381	92.089	61.901	1340.086

Summary of bottom-50 countries' population size:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.0196	0.3520	3.4213	10.2194	9.9751	161.0137

Second parameter we like to analyse in the same matter is the income per capita (in \$1.000):

Summary of all countries' income per capita:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.216	1.526	5.638	15.728	17.401	172.676

Summary of top-50 countries' income per capita:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.360	5.699	12.890	22.677	38.330	97.254

Summary of bottom-50 countries' income per capita:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.2160	0.8662	3.9100	12.0196	14.1977	97.0000

## Combining Population Size and Income Parameters

On our plots we see some countries as obvious exceptions which are over- or underperforming in comparison to their population or income level parameters. Some examples for those countries are:

Averages of all participating countries as a reference:

##	Income	Popsiz
##	Min. : 0.216	Min. : 0.0094
##	1st Qu.: 1.526	1st Qu.: 1.2306
##	Median : 5.638	Median : 6.4980
##	Mean : 15.728	Mean : 34.0245
##	3rd Qu.: 17.401	3rd Qu.: 22.6124
##	Max. : 172.676	Max. : 1340.0865

Pakistan, 7th most populous country in the world:

##	BordaPoints	Popsiz	Income
## Pakistan	0	180.4564	1.201

Nigeria, 8th most populous country in the world:

##	BordaPoints	Popsiz	Income
## Nigeria	0	170.1175	1.49

Monaco, 1st highest income country in the world:

```
##          BordaPoints Popsiz Income
## Monaco          0  0.0364 172.676
```

Liechtenstein, 2nd highest income country in the world:

```
##          BordaPoints Popsiz Income
## Liechtenstein      0  0.0363 143.151
```

Jamaica, 138th most populous country in the world:

```
##          BordaPoints Popsiz Income
## Jamaica          24 2.887784  5.402
```

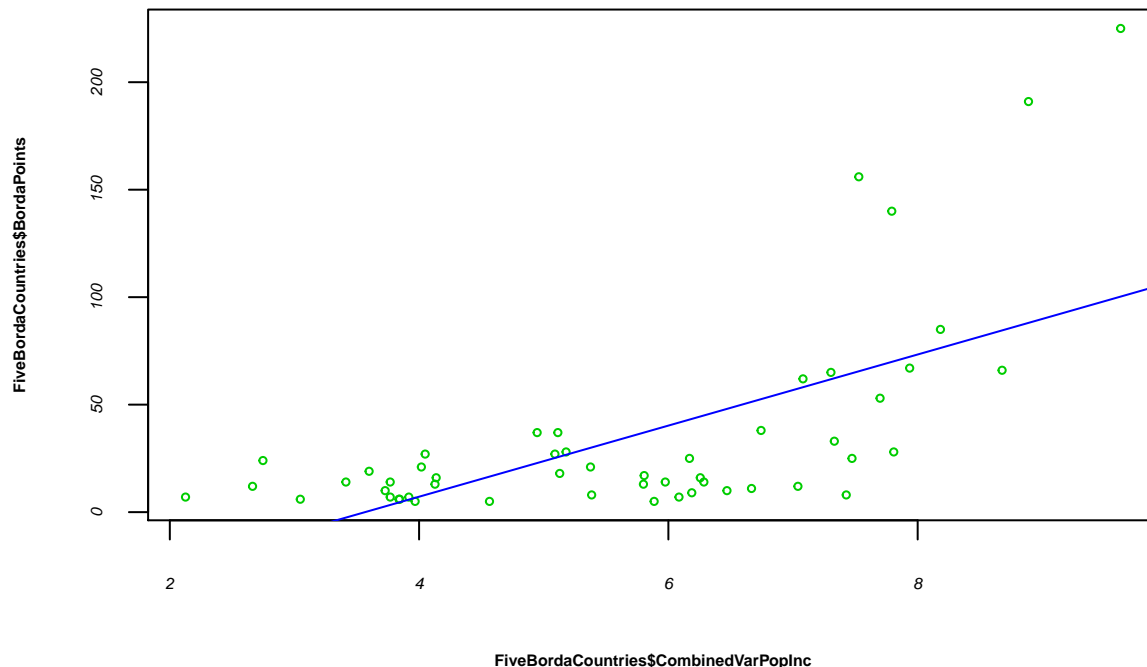
Ethiopia, 6th lowest income country in the world:

```
##          BordaPoints Popsiz Income
## Ethiopia         14 84.30907  0.36
```

Answering the second question of our analysis may help us to combine population size and income parameters with correct weights and coming up with a fitting regression model can help us explain these inconsistencies.

An experimental plot to combine these two parameters and create a better correlation would be:

*Distribution of Borda Points by  $\text{Ln}(\text{PopSize}) + (\text{Ln}(\text{Income}))$  where a country has  $\geq 5$  Borda Points*



## Summary

After our analysis it can be argued that the population size and income per capita parameters of a country have an obvious effect on the countries' success in olympics. However as in many socioeconomic topics only two parameters cannot directly explain an outcome in every case. It's obvious that there are many exceptions or extreme cases where one of the parameters or both of them show a contrary relationship with the success in olympics. Thus it can be stated that there is still enough room for further exploration in this analysis.

## Further Investigation

Our investigation reveals that the data set may possibly be expanded with further parameters about the countries for us to arrive at better conclusions. In this case our further analysis question could be:

**If population size and income are not fully enough to predict the success in olympics, what other measurable metrics can be added to the data to explain the olympics success more precisely?**