

How Institutional & Economic Factors Affect Scientific Output

Omar Hassan
Nichola Millman
Zeeshan Javed

April 17, 2021

Contents

1	Introduction	2
2	Data & Methods	2
2.1	Data Collection	2
2.2	Data Cleaning	2
2.3	Methods	2
3	Results	3
3.1	Scientific Output	3
3.2	Institutional Factors	6
3.3	Economic Factors	7
4	Discussion & Conclusion	8
5	Appendix	9
5.1	Tweet	9
5.2	Datasets & remove.txt	9
5.3	Code	9

1 Introduction

Scientific output has seen a tremendous rise over the past few years, given this exponential growth, we aim to study the economic and institutional factors which affect scientific output across various countries over the years, and we also aim to determine which countries have the greatest scientific output. To do this, we first define scientific output qualitatively and quantitatively, using these two definitions, we explore various economic and institutional aspects which would affect the scientific output by country and across the years. Our economic factors are GDP per capita and research & development expenditure, while our institutional factors are researchers per million and the education index. We then compare these factors with our two functions for scientific output, with the qualitative function being the number of citations per document, whereas the quantitative function is the number of documents per million.

2 Data & Methods

2.1 Data Collection

The majority of our data was collected using an API called WDI. WDI is the API used by the world bank to export all of its data to anyone that wants access. Each data set on the world bank has an indicator. Entering this indicator into the API would return the corresponding data set. An example of an indicator is SP.POP.TOTL; this indicator returns the total population of each country. The API also allowed us to collect data for multiple years at the same time by giving the API a start and an end date. All of the data returned was already formatted into a data frame. We used the world bank API to collect data on population, GDP, Research & development expenditure, Science & technical articles, and researchers per million of each country. We collected all this data over the years of 2012 to 2018. Although the majority of our data was collected from the world bank, we did use other sources to manually collect the citations and education index data. The citations data was collected from SCImago Institutions Rankings and the education index data was collected from the United Nations Development Programme.

2.2 Data Cleaning

Even though collecting the data was relatively simple, cleaning the data was not such an easy job. A substantial amount of the data that we received from the world bank was missing important entries. Typically countries with a lower population were missing the most amount of data, to counter this issue we decided to narrow our data set down to the countries that had a population greater than one million. This did remove a lot of the countries with missing data but there were still some missing entries, our solution was to calculate the missing values using the sample mean of all the previous years where we did have the data. This resolved all of the remaining missing values. Next, we decided to rename all of the columns in the data set to give them more descriptive and usable names. Finally, some of the data points were corresponding to regions of the world or class status rather than countries. So, we removed all regional and class data and kept only the data for the countries.

2.3 Methods

Our most effective tool was data visualizations and correlations. Data visualization often allows us to not only see the relationship between two variables but also understand the limitations of the data (we can limit our data to gain a deeper understanding of the relationship subject to some condition). Our main methods for data visualizations involved scatter, line, bar, and boxplots. Besides data visualizations, we also look at the correlations to determine the strength of the relationships between the variables. When there's an uncertainty regarding the statistical significance of the relationship, we perform a regression analysis (linear regression) and check the confidence intervals to determine the significance.

3 Results

3.1 Scientific Output

We consider two main factors as a measure of scientific output; documents published per million people, and citations per document. The number of documents published per million is a quantitative measure that accounts for the total amount of academic articles being published without being distorted by a given country's population, while the number of citations per article is a qualitative measure. We first see how the number of documents per million change by year:

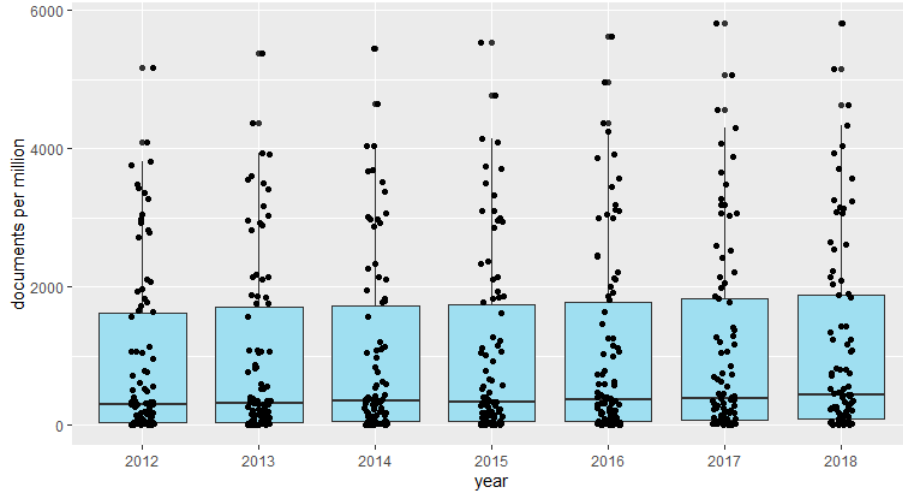


Figure 1: Documents per million by year

As can be seen by figure 1, there appears to have been a slight increase in the third quartile over time, however, after performing a linear regression, it is shown to not be statistically significant.

We now examine how citations per document changes over time:

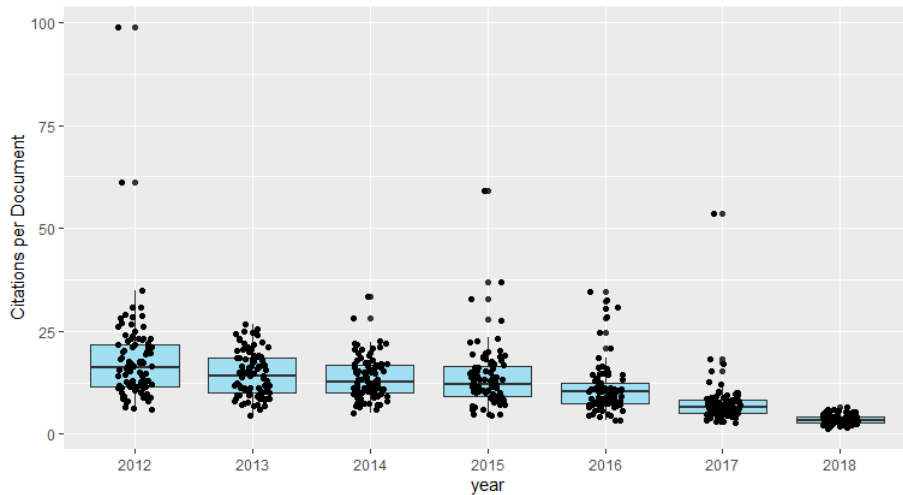


Figure 2: Citations per document by year

Figure 2 shows that there is a clear decrease over the years of the average ratio of citations per document. After performing a linear regression, and having none of the years contain zero in their respective confidence intervals, it is also statistically significant. This suggests that even though the quantity of articles being published around the world has not changed significantly, the amount of

citations they have received has decreased. One possible explanation for this would be that older documents have more time to be cited than newer ones.

We now look at the relationship between documents per million and citations per document; we expect to have a strong correlation:

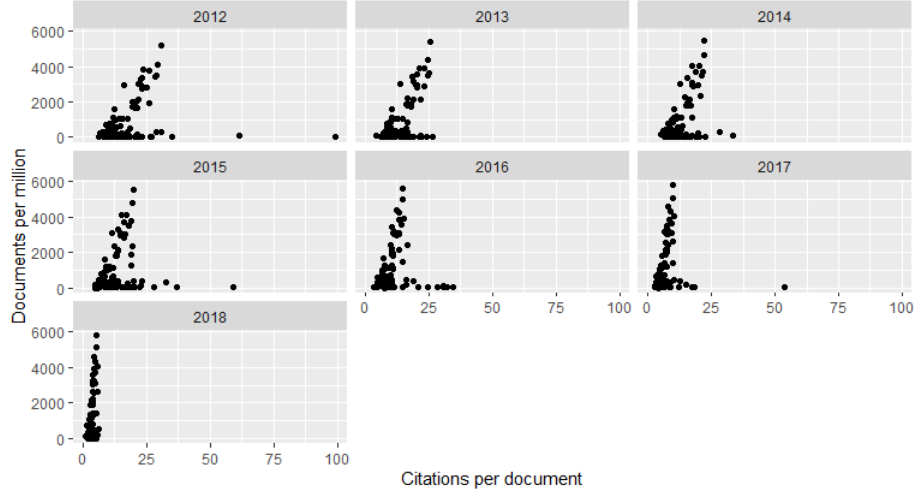


Figure 3: Documents per million vs Citations per document

There appears to be a slightly L shaped relationship between these two variables. With a very low linear correlation value of 0.08713071, this implies that there is no linear relationship. However, one explanation for this would be that a few well cited articles from a country which produces relatively few documents each year would have a large effect on the country's value of citations per document. When considering the countries that have an above average number of documents per million, there appears to be a much stronger correlation:

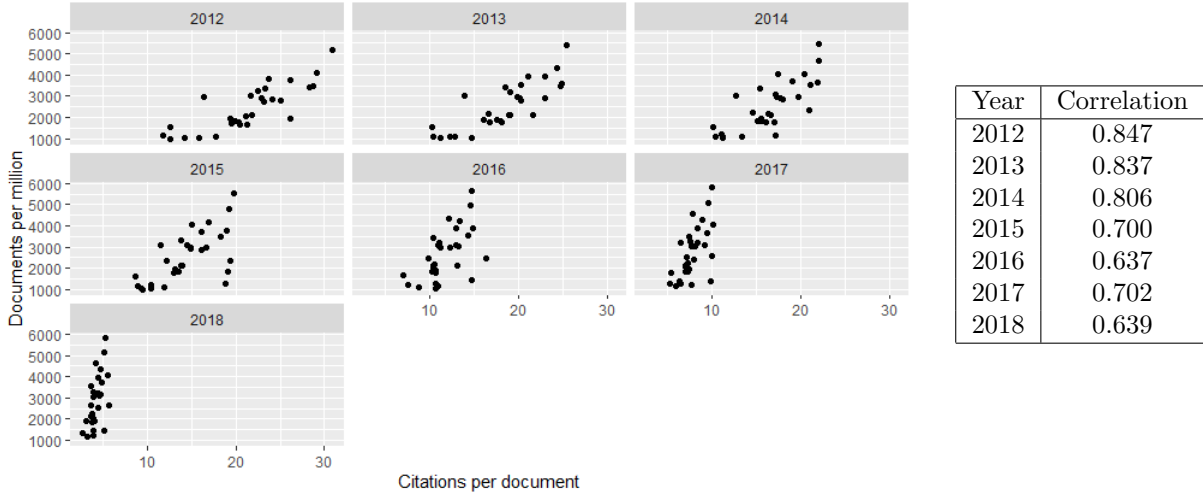


Figure 4: Documents per million (above average) vs Citations per document

As can be seen by figure 4, there appears to be a solid correlation between documents per million and citations per document among the countries which produced an above average amount of documents per million. It is also notable that the slope increases over time, which can be attributed to the decreasing value of citations per document over time.

We now examine which countries have the greatest scientific output both qualitatively and quantitatively; we only plot the top 10 countries for clarity:

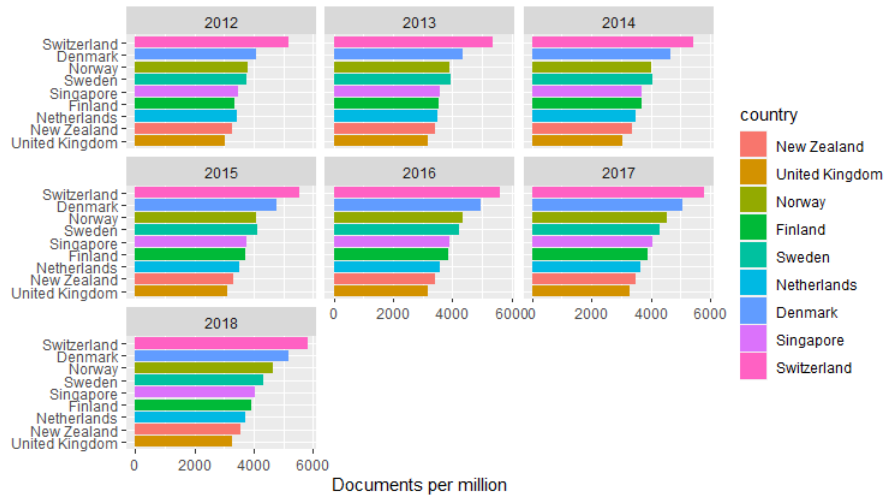


Figure 5: Top 10 Countries by Documents per million by year

It is clear that Switzerland is the top performer across all the years. Followed by Denmark and then Norway/Sweden. There seems to be a greater presence of Nordic countries compared to any other region.

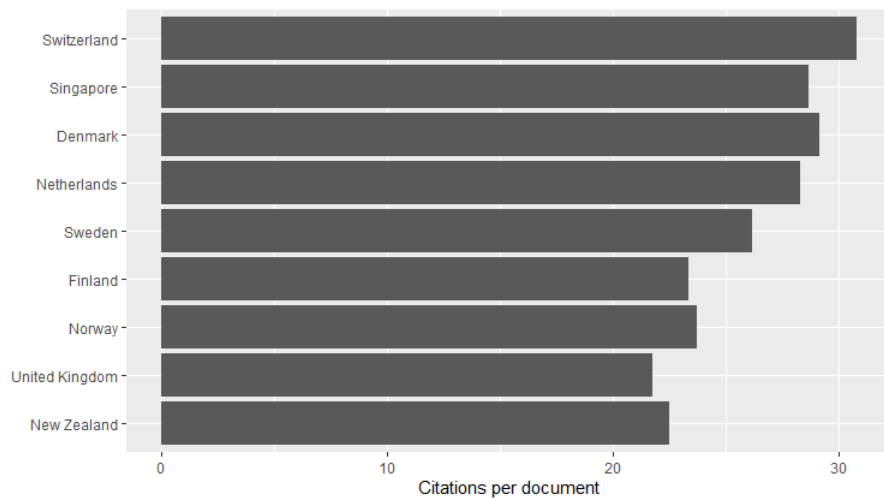


Figure 6: Top 10 Countries by Total Citations per document across 2012-2018

Switzerland is the top performer when it comes to the total citations per document across the years, and the Nordic regions maintain a strong presence in the top 10. We can confidently say that Switzerland is the country with the greatest scientific output, while the Nordic region has the most scientific output when compared to other regions. We now look to the institutional and economic factors which may explain why these countries have such a high scientific output.

3.2 Institutional Factors

We explore two main relationships, the relationship between documents per million and both researchers per million and the education index. Starting with the former, we perform a simple scatter plot for each year giving us:

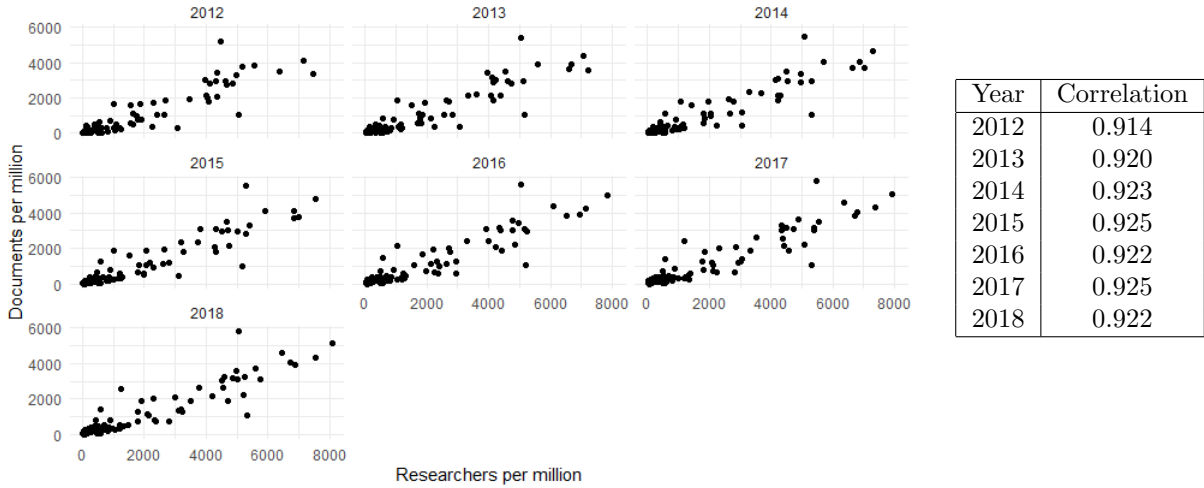


Figure 7: Documents per million vs Researchers per million

As can be expected, there is a strong correlation between the number of documents published and the number of researchers per million. This is visible from figure 7, the correlation confirms this relationship and can be found to be on average about 0.92. This correlation remains fairly constant throughout the years.

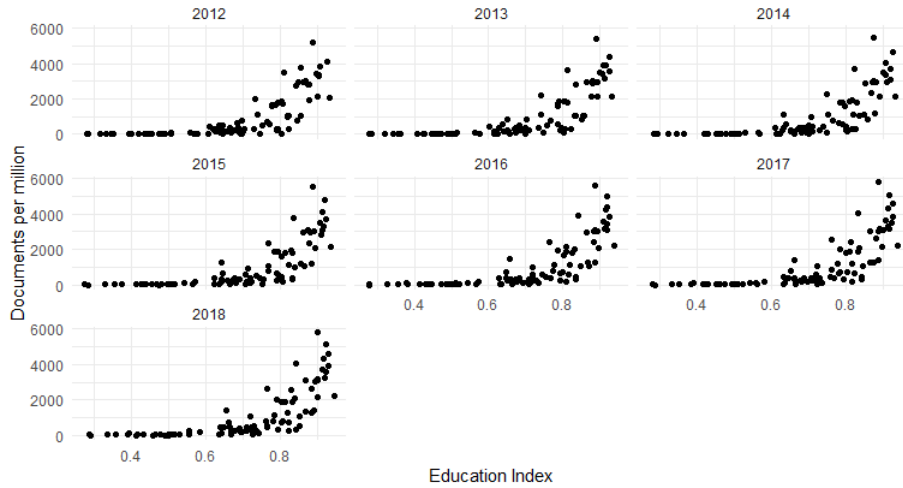


Figure 8: Documents per million vs Education Index

We now explore the relationship between documents per million and education index, we expect countries with a greater education index to produce more documents per million. The relationship is given in figure 8. It is quite clear that there is no relation up until the education index is around 0.6, to get a more accurate investigation of this relationship, we can filter all data points with an education index of less than 0.6. This way, we can get a meaningful correlation between education index greater than 0.6 and documents per million.

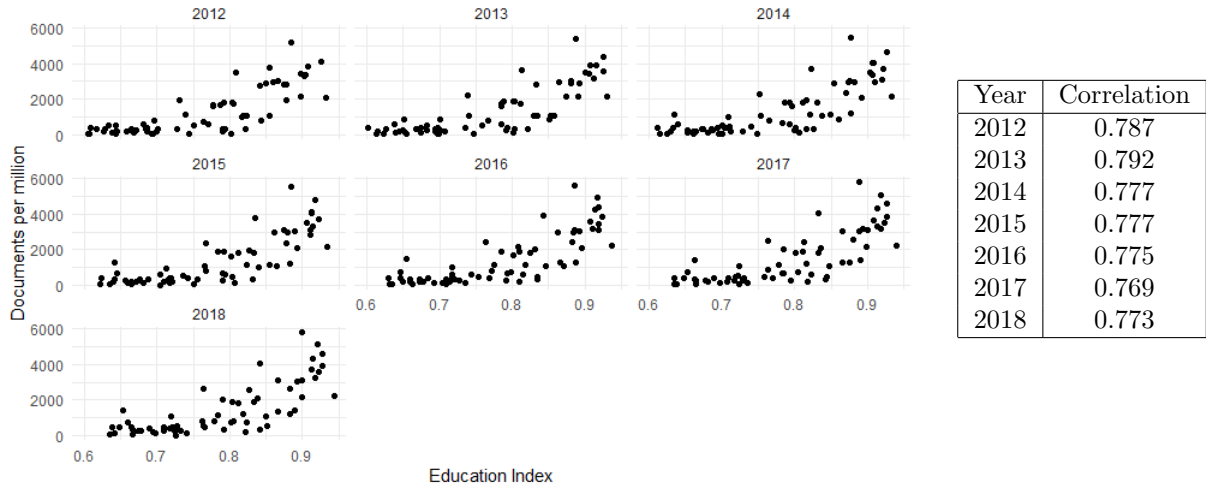


Figure 9: Documents per million vs Education Index (> 0.6)

After filtering, we can see that the relationship between the two variables is quite strong, the exponential relation appears to be more linear when focusing on education index values of > 0.6 . The correlation further proves the strength of this relationship with a somewhat constant value across the years, averaging out at around 0.77.

3.3 Economic Factors

To measure the effects of economic factors on scientific output, we decided to measure the effects of 2 factors on scientific output, GDP per capita and research & development expenditure. Starting with the former:

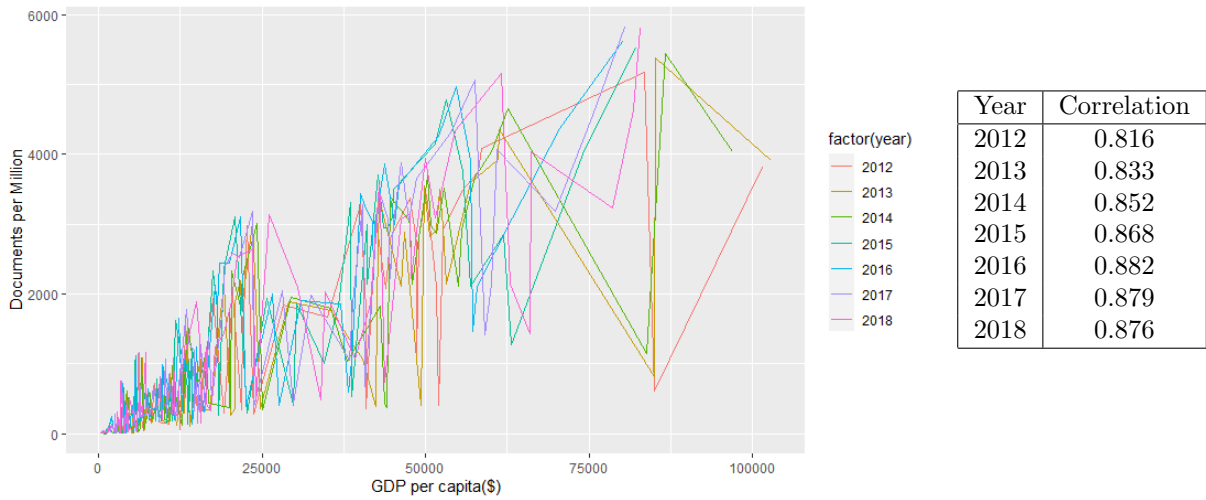


Figure 10: Documents per million vs GDP per capita

As we can clearly see from the above figure, as GDP per capita of a country increases, the amount of scientific output also increases. There seems to be a very strong correlation between GDP per capita and Scientific output. The correlation also seems to be increasing over the years.

We now look to research & development expenditure when compared to documents per million, we expect to have quite a strong correlation given that the more a country spends on research and development, the more documents it should produce:

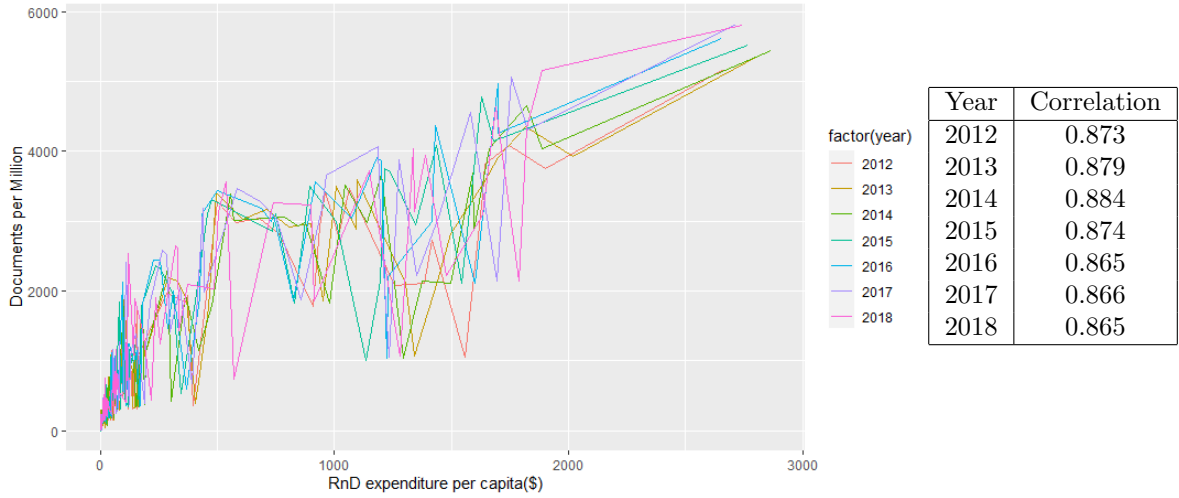


Figure 11: Documents per million vs Research & Development Expenditure

We can also see there's a strong correlation from figure 11. The correlation between the years has stayed relatively the same. One possible explanation could be that the research & development expenditure did not change significantly through the years of 2012 to 2018.

4 Discussion & Conclusion

After having performed the analysis, we have found a positive correlation between the number of documents published each year and the number of citations per document for countries which publish documents above the average mean. We also discovered that the country with the largest documents per million was consistently Switzerland from the years 2012-2018. In addition, we discovered a positive correlation between GDP per capita and documents per million, and an exponential relationship between the education index and documents per million. The strongest correlation we found was between researchers per million and documents per million.

Some of the limitations that we experienced were the lack of available data for many of the countries in the world. In order to offset this, some of the countries with a few missing data points used the mean average over the years to miss in filling values, but some countries were missing values altogether. Additionally, the numbers provided in the SCImango dataset for citations and number of documents published each year are most likely an underestimation as they were compiled through a collection of academic journals, which was most likely not comprehensive.

Throughout the analysis, there were both systematic and random sources of error. By filling in the missing values of some of the columns with the mean average value across the years, the impact of outliers in the calculations of correlation and linear regression may have been underestimated, in particular for year by year comparisons. Additionally, some academic articles and citations being excluded in the scimango dataset could be a source of random error, as, for the sake of comparison, the impact of the articles and citations that were not included could either increase or decrease a given country's value for citations per document.

We would like to explore this data further using better methods of correcting our missing data, we also lost a significant amount of countries in the cleaning process. We could potentially do fitted values rather than the sample mean and compare the results with this report, and coupling this with finding more comprehensive datasets could provide us with the data to explore our questions to a

greater extent. It would also be worth exploring scientific output by region rather than by country, we saw that the Nordic region seemed to be over-represented in the high output countries, we could explore this idea with various other regions as well.

In conclusion, we have found that researchers per million, GDP, research spending, and the education index of a given country all influence its scientific output, with both the quantity and the quality of the countries with a larger number of documents per million being closely related.

5 Appendix

5.1 Tweet

Why do Switzerland & Denmark have the highest scientific output globally? We explored the institutional and economic factors which lead to a high scientific output, and the prevailing factors are the number of researchers, the GDP per capita, and the research expenditure.

5.2 Datasets & remove.txt

The datasets and the remove.txt file can be found on: <https://github.com/S-tuberosum/Scientific-Output>

5.3 Code

```
1  # Code written by Zeeshan Javed, Nichola Millman, and Omar Hassan
2  #-----
3  #           READ ME
4  #-----
5  #Before you do anything, make sure you have all of the libraries installed
6  #Also make sure all the provided .xlsx files are in the correct directory
7  #After doing the above, run the entire script at once and all the data sets will be created
8
9  #We are only using data for countries with a population greater than 1 million
10
11 #There is missing data for some countries
12
13 #Run any of your experiments at the very bottom of the script
14 #so that it doesn't mess with the created datasets
15
16 #The data frame labeled 'FinalData' contains all the data for all the years
17
18 #The data frame labeled 'FinalDataYYYY' contains all the data for the year YYYY
19
20 library(WDI)
21 library(tidyverse)
22 library(dplyr)
23 library(readxl)
24 library(dslabs)
25 library(broom)
26 library(cowplot)
27
28
29 #-----
30 #           GET THE DATA
31 #-----
32
33
34 #population of each country from 2012 to 2018
```

```

35 population = WDI(indicator='SP.POP.TOTL', start=2012, end=2018)
36
37 #GDP of each country from 2012 to 2018
38 gdp = WDI(indicator='NY.GDP.MKTP.CD', start=2012, end=2018)
39
40 #research and development expenditure (% of GDP) from 2012 to 2018
41 RnD_expenditure = WDI(indicator='GB.XPD.RSDV.GD.ZS', start=2012, end=2018)
42
43 #Scientific and technical journal articles from 2012 to 2018
44 Sci_and_tech_journals = WDI(indicator = 'IP.JRN.ARTC.SC', start=2012, end=2018)
45
46 #researchers per million from 2012 to 2018
47 researchers_million = WDI(indicator = 'SP.POP.SCIE.RD.P6', start =2012, end =2018)
48
49 #Citations per country / add the year column to the data
50 Citations2012 <- read_excel("./Data/scimagojr country rank 2012.xlsx")
51 Citations2012$year = 2012
52
53 Citations2013 <- read_excel("./Data/scimagojr country rank 2013.xlsx")
54 Citations2013$year = 2013
55
56 Citations2014 <- read_excel("./Data/scimagojr country rank 2014.xlsx")
57 Citations2014$year = 2014
58
59 Citations2015 <- read_excel("./Data/scimagojr country rank 2015.xlsx")
60 Citations2015$year = 2015
61
62 Citations2016 <- read_excel("./Data/scimagojr country rank 2016.xlsx")
63 Citations2016$year = 2016
64
65 Citations2017 <- read_excel("./Data/scimagojr country rank 2017.xlsx")
66 Citations2017$year = 2017
67
68 Citations2018 <- read_excel("./Data/scimagojr country rank 2018.xlsx")
69 Citations2018$year = 2018
70
71 #Education index per country / add the year column to the data, change col names
72
73 EducationIndex2012 <- read_excel("./Data/EducationIndex2012.xlsx")
74 EducationIndex2012$year = 2012
75
76 EducationIndex2013 <- read_excel("./Data/EducationIndex2013.xlsx")
77 EducationIndex2013$year = 2013
78
79 EducationIndex2014 <- read_excel("./Data/EducationIndex2014.xlsx")
80 EducationIndex2014$year = 2014
81
82 EducationIndex2015 <- read_excel("./Data/EducationIndex2015.xlsx")
83 EducationIndex2015$year = 2015
84
85 EducationIndex2016 <- read_excel("./Data/EducationIndex2016.xlsx")
86 EducationIndex2016$year = 2016
87
88 EducationIndex2017 <- read_excel("./Data/EducationIndex2017.xlsx")
89 EducationIndex2017$year = 2017
90
91 EducationIndex2018 <- read_excel("./Data/EducationIndex2018.xlsx")

```

```

92 EducationIndex2018$year = 2018
93 names(EducationIndex2018)[names(EducationIndex2018) == 'education_vaue'] <- 'education_value'
94
95 #now bind all the Citations and EducationIndex into 1 data set
96 citations <- rbind(Citations2012,Citations2013,Citations2014,Citations2015,Citations2016,
97                   Citations2017,Citations2018)
98
99 eduindex <- rbind(EducationIndex2012,EducationIndex2013,EducationIndex2014,
100                  EducationIndex2015,EducationIndex2016,EducationIndex2017,
101                  EducationIndex2018)
102
103 #now delete some columns from the Citations & eduindex and rename some columns
104 drop <- c('Rank','Region')
105 citations = citations[!(names(citations) %in% drop)]
106 names(citations)[1] <- 'country'
107
108 names(eduindex)[names(eduindex) == 'Country'] <- 'country'
109
110 #-----
111 #           Create the Super data set
112 #-----
113
114
115 FinalData <- inner_join(population, gdp, by=c('country','year','iso2c'))
116 FinalData <- inner_join(FinalData, RnD_expenditure, by=c('country','year','iso2c'))
117 FinalData <- inner_join(FinalData, Sci_and_tech_journals, by=c('country','year','iso2c'))
118 FinalData <- inner_join(FinalData, researchers_million, by=c('country','year','iso2c'))
119
120 #rename some columns then add the citation data
121 names(FinalData)[1] <- 'symbol'
122 names(FinalData)[3] <- 'population'
123 names(FinalData)[5] <- 'gdp'
124 names(FinalData)[6] <- 'RnD_Expenditure'
125 names(FinalData)[7] <- 'sci_tech_articles'
126 names(FinalData)[8] <- 'researchers_per_million'
127
128 #now add the citations data
129 FinalData <- inner_join(FinalData, citations, by=c('country','year'))
130 FinalData <- inner_join(FinalData, eduindex, by=c('country','year'))
131
132 names(FinalData)[13] <- 'cit_per_doc'
133 names(FinalData)[15] <- 'education_index'
134
135 #-----
136 #           Clean the Super data set
137 #-----
138
139
140 #now remove all countries with a population less than 1 million
141 FinalData <- subset(FinalData, population > 1000000)
142
143 #create data frame of countries with NA across all the years and remove them
144 remove <- read.delim("remove.txt")
145 FinalData <- anti_join(FinalData, remove, by='country')
146
147
148 #replacing all NA values with the sample mean of the column

```

```

149 FinalData <- group_by(FinalData, country)
150 FinalData <- mutate(FinalData, mean_rnd = mean(RnD_Expenditure, na.rm = TRUE),
151                    mean_rpm = mean(researchers_per_million, na.rm = TRUE))
152
153 FinalData$RnD_Expenditure <- ifelse(is.na(FinalData$RnD_Expenditure),
154                                    FinalData$mean_rnd, FinalData$RnD_Expenditure)
155
156 FinalData$researchers_per_million <- ifelse(is.na(FinalData$researchers_per_million),
157                                            FinalData$mean_rpm, FinalData$researchers_per_million)
158
159 #removing columns we don't need & formatting
160 FinalData <- FinalData[-c(1, 14, 16:17)]
161 FinalData$education_index <- as.double(FinalData$education_index)
162
163 #adding our scientific output function
164
165 FinalData <- FinalData %>%
166   mutate(doc_per_mil = (Documents*1000000)/population)
167
168 #add gdp_per_capita to the Final data
169 FinalData = FinalData%>%
170   mutate(gdp_per_capita = (gdp/population))
171
172 #add RnD expenditure dollar amount
173 FinalData = FinalData%>%
174   mutate(RnD_dollar_amount = (RnD_Expenditure/100) * gdp)
175
176 #add RnD expenditure per capita to the FinalData
177 FinalData = FinalData%>%
178   mutate(RnD_per_capita = (RnD_dollar_amount/population))
179
180 #-----
181 # Delete all the 'extra' data sets that we don't need
182 #-----
183
184
185 rm(Citations2012)
186 rm(Citations2013)
187 rm(Citations2014)
188 rm(Citations2015)
189 rm(Citations2016)
190 rm(Citations2017)
191 rm(Citations2018)
192 rm(citations)
193 rm(EducationIndex2012)
194 rm(EducationIndex2013)
195 rm(EducationIndex2014)
196 rm(EducationIndex2015)
197 rm(EducationIndex2016)
198 rm(EducationIndex2017)
199 rm(EducationIndex2018)
200 rm(eduindex)
201 rm(gdp)
202 rm(drop)
203 rm(population)
204 rm(RnD_expenditure)
205 rm(Sci_and_tech_journals)

```

```

206 rm(researchers_million)
207
208
209 #-----
210 # Create the separate data sets for each year
211 #-----
212
213
214 FinalData2012 <- FinalData %>%
215   filter(year == 2012)
216
217 FinalData2013 <- FinalData %>%
218   filter(year == 2013)
219
220 FinalData2014 <- FinalData %>%
221   filter(year == 2014)
222
223 FinalData2015 <- FinalData %>%
224   filter(year == 2015)
225
226 FinalData2016 <- FinalData %>%
227   filter(year == 2016)
228
229 FinalData2017 <- FinalData %>%
230   filter(year == 2017)
231
232 FinalData2018 <- FinalData %>%
233   filter(year == 2018)
234
235
236 #-----
237 #                               analyses
238 #-----
239
240 #-----
241 #       Scientific Output
242 #-----
243
244 #First lets take a look at the relationship between number of documents published per million per o
245
246 output_over_time = FinalData%>%
247   group_by(year)%>%
248   summarize(sum(Documents),mean(Documents))
249
250 #boxplot no jitter
251 FinalData$year = factor(FinalData$year)
252 ggplot(FinalData,aes(x= year, y = doc_per_mil)) +
253   geom_boxplot( fill = "#9FDFF1") +
254   ylab("Documents per million")
255
256 #boxplot with jitter
257 FinalData$year = factor(FinalData$year)
258 ggplot(FinalData,aes(x= year, y = doc_per_mil)) +
259   geom_boxplot( fill = "#9FDFF1") +
260   geom_jitter(height =0.10,width =0.10) +
261   ylab("documents per million")
262

```

```

263 #linear model
264 fit <- lm(doc_per_mil~year, data = FinalData)
265 fit
266 confint(fit)
267
268
269 #Now, let's look at the relationship between citations and articles published
270
271
272 #boxplot not jitter
273 FinalData$year = factor(FinalData$year)
274 ggplot(FinalData,aes(x= year, y = cit_per_doc)) +
275   geom_boxplot( fill = "#9FDF1" ) +
276   ylab("Citations per Document")
277
278 #boxplot with jitter
279 FinalData$year = factor(FinalData$year)
280 ggplot(FinalData, aes(x = year, y = cit_per_doc)) +
281   geom_boxplot(fill = "#9FDF1" ) +
282   geom_jitter(height =0.15, width =0.15) +
283   ylab("Citations per document")
284
285 #linear model
286 fit2 <- lm(cit_per_doc~year, data = FinalData)
287 fit2
288 confint(fit2)
289
290 #Now let's look at output vs population
291
292 #point graph over the years
293
294 FinalData%>%
295   ggplot()+
296   geom_point(aes(x= population/100000, y = Documents/1000))+
297   xlab("Population")+
298   ggtitle("Population vs Published Documents")+
299   facet_wrap(~year)+
300   theme(plot.title = element_text(hjust = 0.5))
301
302 #point graph zoomed in (outliers out of view)
303 FinalData2018%>%
304   ggplot()+
305   geom_point(aes(x= population/1000000, y = Documents/100))+
306   coord_cartesian(ylim = c(0,3000))+
307   ggtitle("Population vs Published Documents 2018")+
308   theme(plot.title = element_text(hjust = 0.5))
309
310 #linear model and correlation
311 cor(FinalData$population,FinalData$Documents,method = "pearson")
312 fit3 <- lm(population~Documents, data = FinalData)
313 fit3
314 confint(fit3)
315
316 #Now lets look at the top performers
317
318 #first, lets graph citations per document vs documents per million
319 FinalData%>%

```

```

320   ggplot()+
321   geom_point(aes(x= cit_per_doc, y = doc_per_mil))+
322   xlab("Citations per document")+
323   ylab("Documents per million")+
324   facet_wrap(~year)
325
326   cor(FinalData$cit_per_doc,FinalData$doc_per_mil,method = "pearson")
327   fit4 <- lm(cit_per_doc~doc_per_mil, data = FinalData)
328   fit4
329   confint(fit4)
330
331   #Now lets look at above average countries
332
333   #above mean average docs per mil dataset
334   above_avg2012 = FinalData2012%>%
335     filter(doc_per_mil>mean(FinalData2012$doc_per_mil))
336
337   above_avg2013 = FinalData2013%>%
338     filter(doc_per_mil>mean(FinalData2013$doc_per_mil))
339
340   above_avg2014 = FinalData2014%>%
341     filter(doc_per_mil>mean(FinalData2014$doc_per_mil))
342
343   above_avg2015 = FinalData2015%>%
344     filter(doc_per_mil>mean(FinalData2015$doc_per_mil))
345
346   above_avg2016 = FinalData2016%>%
347     filter(doc_per_mil>mean(FinalData2016$doc_per_mil))
348
349   above_avg2017 = FinalData2017%>%
350     filter(doc_per_mil>mean(FinalData2017$doc_per_mil))
351
352   above_avg2018 = FinalData2018%>%
353     filter(doc_per_mil>mean(FinalData2018$doc_per_mil))
354
355   above_avg = full_join(above_avg2012,above_avg2013)%>%
356     full_join(above_avg2014)%>%
357     full_join(above_avg2015)%>%
358     full_join(above_avg2016)%>%
359     full_join(above_avg2017)%>%
360     full_join(above_avg2018)
361
362   #graph of above average docs per mil vs citations per document
363
364   above_avg%>%
365     ggplot()+
366     geom_point(aes(x= cit_per_doc, y = doc_per_mil))+
367     xlab("Citations per document")+
368     ylab("Documents per million")+
369     facet_wrap(~year)
370
371   above_avg_cor = above_avg%>%
372     group_by(year)%>%
373     summarize(cor(cit_per_doc,doc_per_mil,method = "pearson"))
374
375   #performing a linear regression to see if the above mean average docs_per_mil are increasing over t
376   fit4 <- lm(year~doc_per_mil, data = above_avg)

```



```

377 fit4
378 confint(fit4)
379
380 #Now, lets consider the top ten countries for docs per million over the years
381 #first let's filter the dataset for the top ten each year
382 top_performers2012 = FinalData2012%>%
383   filter(doc_per_mil>3000)
384
385 top_performers2013 = FinalData2013%>%
386   filter(doc_per_mil>3100)
387
388 top_performers2014 = FinalData2014%>%
389   filter(doc_per_mil>3050)
390
391 top_performers2015 = FinalData2015%>%
392   filter(doc_per_mil>3104.5)
393
394 top_performers2016 = FinalData2016%>%
395   filter(doc_per_mil>3150)
396
397 top_performers2017 = FinalData2017%>%
398   filter(doc_per_mil>3200)
399
400 top_performers2018 = FinalData2018%>%
401   filter(doc_per_mil>3250)
402
403
404 top_performers = full_join(top_performers2012,top_performers2013)%>%
405   full_join(top_performers2014) %>%
406   full_join(top_performers2015) %>%
407   full_join(top_performers2016) %>%
408   full_join(top_performers2017) %>%
409   full_join(top_performers2018)
410
411 #now let's do a bar graph of each year
412
413 top_performers%>%
414   ggplot(aes(y= doc_per_mil,x= reorder(country,doc_per_mil), fill = country))+
415   geom_bar(stat = "identity",position = "dodge")+
416   ylab("Documents per million")+
417   xlab("")+
418   coord_flip() +
419   facet_wrap(~year)
420
421 #Now lets look at citations per doc of the top ten countries
422
423 top_performers %>%
424   ggplot(aes(y= cit_per_doc,x = country)) +
425   geom_bar(stat = "identity",position = "dodge") +
426   ylab("Citations per document") +
427   xlab("") +
428   coord_flip()
429
430
431 #-----
432 #           Institutional factors vs scientific output
433 #-----

```

```

434
435
436 # scatter doc/mil vs researchers/mil
437 FinalData %>%
438   ggplot() +
439   geom_point(aes(x= researchers_per_million, y = doc_per_mil)) +
440   xlab("Researchers per million") +
441   ylab("Documents per million") +
442   facet_wrap(~year) +
443   theme_minimal()
444
445 # correlations for doc/mil vs researchers/mil
446 temp1 <- cor.test(FinalData2012$researchers_per_million,
447                   FinalData2012$doc_per_mil, method = 'pearson')
448 temp1
449
450 temp2 <- cor.test(FinalData2013$researchers_per_million,
451                   FinalData2013$doc_per_mil, method = 'pearson')
452 temp2
453
454 temp3 <- cor.test(FinalData2014$researchers_per_million,
455                   FinalData2014$doc_per_mil, method = 'pearson')
456 temp3
457
458 temp4 <- cor.test(FinalData2015$researchers_per_million,
459                   FinalData2015$doc_per_mil, method = 'pearson')
460 temp4
461
462 temp5 <- cor.test(FinalData2016$researchers_per_million,
463                   FinalData2016$doc_per_mil, method = 'pearson')
464 temp5
465
466 temp6 <- cor.test(FinalData2017$researchers_per_million,
467                   FinalData2017$doc_per_mil, method = 'pearson')
468 temp6
469
470 temp7 <- cor.test(FinalData2018$researchers_per_million,
471                   FinalData2018$doc_per_mil, method = 'pearson')
472 temp7
473
474 # scatter doc/mil vs education index
475 FinalData %>%
476   filter(education_index > 0.6) %>%
477   ggplot() +
478   geom_point(aes(x= education_index, y = doc_per_mil)) +
479   xlab("Education Index") +
480   ylab("Documents per million") +
481   facet_wrap(~year) +
482   theme_minimal()
483
484 # correlations for doc/mil vs education index
485 FinalData %>%
486   filter(education_index > 0.6) %>%
487   group_by(year)
488
489 temp8 <- cor.test(filter(FinalData2012, education_index > 0.6)$education_index,
490                   filter(FinalData2012, education_index > 0.6)$doc_per_mil,

```

```

491         method = 'pearson')
492 temp8
493
494 temp9 <- cor.test(filter(FinalData2013, education_index >0.6)$education_index,
495                   filter(FinalData2013, education_index >0.6)$doc_per_mil,
496                   method = 'pearson')
497 temp9
498
499 temp10 <- cor.test(filter(FinalData2014, education_index >0.6)$education_index,
500                   filter(FinalData2014, education_index >0.6)$doc_per_mil,
501                   method = 'pearson')
502 temp10
503
504 temp11 <- cor.test(filter(FinalData2015, education_index >0.6)$education_index,
505                   filter(FinalData2015, education_index >0.6)$doc_per_mil,
506                   method = 'pearson')
507 temp11
508
509 temp12 <- cor.test(filter(FinalData2016, education_index >0.6)$education_index,
510                   filter(FinalData2016, education_index >0.6)$doc_per_mil,
511                   method = 'pearson')
512 temp12
513
514 temp13 <- cor.test(filter(FinalData2017, education_index >0.6)$education_index,
515                   filter(FinalData2017, education_index >0.6)$doc_per_mil,
516                   method = 'pearson')
517 temp13
518
519 temp14 <- cor.test(filter(FinalData2018, education_index >0.6)$education_index,
520                   filter(FinalData2018, education_index >0.6)$doc_per_mil,
521                   method = 'pearson')
522 temp14
523 #-----
524 #      Economic factors vs scientific output
525 #-----
526
527
528 #graph gdp per capita vs documents per million over the years
529 FinalData %>%
530   ggplot(aes(x = gdp_per_capita, y= doc_per_mil,group=year, color = factor(year))) +
531   xlab("GDP per capita($") +
532   ylab("Documents per Million") +
533   geom_line()
534
535 #Find the correlation
536 temp15 <- cor.test(FinalData2012$gdp_per_capita, FinalData2012$doc_per_mil, method = 'pearson')
537 temp15
538
539 temp16 <- cor.test(FinalData2013$gdp_per_capita, FinalData2013$doc_per_mil, method = 'pearson')
540 temp16
541
542 temp17 <- cor.test(FinalData2014$gdp_per_capita, FinalData2014$doc_per_mil, method = 'pearson')
543 temp17
544
545 temp18 <- cor.test(FinalData2015$gdp_per_capita, FinalData2015$doc_per_mil, method = 'pearson')
546 temp18
547

```

```

548 temp19 <- cor.test(FinalData2016$gdp_per_capita, FinalData2016$doc_per_mil, method = 'pearson')
549 temp19
550
551 temp20 <- cor.test(FinalData2017$gdp_per_capita, FinalData2017$doc_per_mil, method = 'pearson')
552 temp20
553
554 temp21 <- cor.test(FinalData2018$gdp_per_capita, FinalData2018$doc_per_mil, method = 'pearson')
555 temp21
556
557 #graph RnD expenditure per capita vs documents per million over the years
558 FinalData %>%
559   ggplot(aes(x=RnD_per_capita,y=doc_per_mil,group=year, color = factor(year))) +
560   xlab("RnD expenditure per capita($)") +
561   ylab("Documents per Million") +
562   geom_line()
563
564
565 #Find the correlation
566 temp22 <- cor.test(FinalData2012$RnD_per_capita, FinalData2012$doc_per_mil, method = 'pearson')
567 temp22
568
569 temp23 <- cor.test(FinalData2013$RnD_per_capita, FinalData2013$doc_per_mil, method = 'pearson')
570 temp23
571
572 temp24 <- cor.test(FinalData2014$RnD_per_capita, FinalData2014$doc_per_mil, method = 'pearson')
573 temp24
574
575 temp25 <- cor.test(FinalData2015$RnD_per_capita, FinalData2015$doc_per_mil, method = 'pearson')
576 temp25
577
578 temp26 <- cor.test(FinalData2016$RnD_per_capita, FinalData2016$doc_per_mil, method = 'pearson')
579 temp26
580
581 temp27 <- cor.test(FinalData2017$RnD_per_capita, FinalData2017$doc_per_mil, method = 'pearson')
582 temp27
583
584 temp28 <- cor.test(FinalData2018$RnD_per_capita, FinalData2018$doc_per_mil, method = 'pearson')
585 temp28

```