

Statistical Data Analysis

Assignment 1

Matt Kuodis & James Zoryk.

2640428 : 2663347

Group: 36

Due Date: 15.02.2022.

Question 1.5.

For the first assignment, we plotted samples from the binomial and Poisson distributions, and compared the resulting histograms. We varied the input parameters for the binomial distribution functions, while keeping $\lambda = \text{size} \times \text{probability}$ as a constant, namely $\lambda = 7.5$. We looked at the histograms with n values at 150 and 500, and for each case we varied the *size* and *probability* inputs, under the restraint of their product being equal to 7.5. We plotted the graphs individually, as plotting multiple histograms would have required us to write a program with built in constants, or with too many inputs.

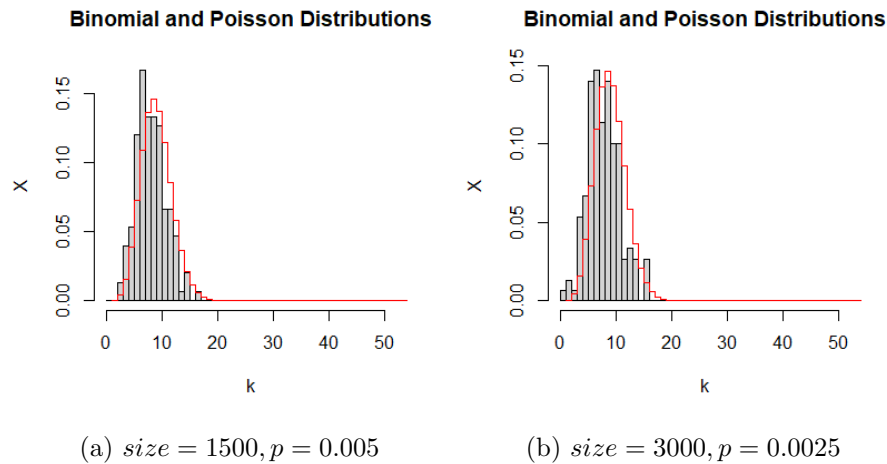


Figure 1: Histograms for $n = 150$, Poisson shown in red

We note that the binomial samples vary quite noticeably under the changes in size and probability, but show no convergence to the Poisson distribution as the *size* variable is increased for small n , specifically $n = 150$. In fact, it seems that the binomial distribution maintains a structure much more similar to that of the Poisson for smaller

size in the case that n is relatively small, but for both histograms we note unexpected peaks and dips that remove any notion of a uniform structure.

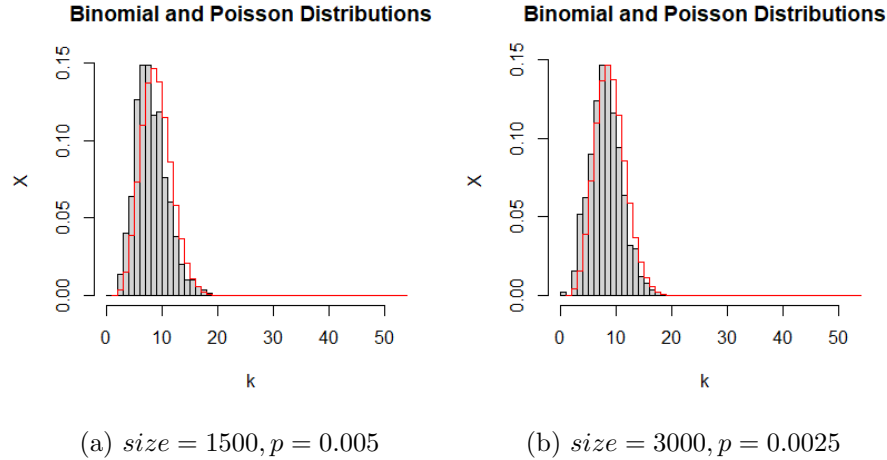


Figure 2: Histograms for $n = 500$, Poisson shown in red

The opposite, however, seems to be true for large n , as it seems that the distributions are much more similar overall, but in particular for larger *size*. The shape of the distribution for $n = 500$ is much more uniformly that of the Poisson, despite that there seems to be a slight deviation to the left in the case of the binomial. An increase in *size* seems to improve the overlap of the two histograms in the latter case, as the left-offset decreases as *size* increases.

Overall, it would seem that the binomial distribution is less unpredictable for large n , and seems to converge to the Poisson as n increases. The *size* hardly seems to alter the likeliness of the two distributions in either case, but appears to reduce the shift slightly to the left for the binomial distribution.

Question 1.6.

In this section we shall make a data summary that will compare the recent partly vaccinated percentage against COVID-19 to that of the Gross Domestic Product (GDP) per capita of countries that are in Asia. We are supplied with up-to-date data from the online source `ourworldindata` from which we carry out our analysis with the assistance of the software package RStudio. The RStudio code that we used to generate the data in this report can be found in the appendix.

Our first step to analysis the data is to process the data. We observe that the data is formatted into a CSV (Comma Separated Values) file, which once imported, the data is organised into named columns. Here we note that the column labeled 'Continent' has the data type of nominal scale, hence it can be grouped. This data type can be use to reduce the data down to countries that lie in Asia, since the raw data contains information about countries world wide, which are not in the interest for this report. Moreover, we are then interested in analysing the data in the two columns labeled 'gdp per capita' and 'partly vacc'. Observing that both these variables are continuous in nature, with 'partly vacc' being a percentage so it is scaled in the range of $[0, 100]$; and 'gdp per capita' being a qualitative ordinal variable. There are a couple of cases where there is no variable give for that specified field and instead it is marked 'NA', we can remove these from our summary while making note about how many we have encounter.

The data set of the two variables 'partly vacc' and 'gdp per capita' can now be numerical analysis with the results located in Table 1 and 2 respectfully.

Sample Size:	61	Sample Size:	61
Mean:	54.67	Mean:	24719
SD:	24.33	SD:	28064.34
Var:	592.32	Var:	787606970
Min:	1.13	Min:	1479
1st qu.:	36.94	1st qu.:	6289
Median:	62.70	Median:	12604
3rd qu.:	73.14	3rd qu.:	35237
Max:	93.45	Max:	116936
IQR:	36.20	IQR:	28947.69
NA's:	17	NA's:	15

Table 1: Numerical analysis for partly vaccinated people in Asia

Table 2: Numerical analysis for GDP per capita in Asia

The tables give us a lot of information about both the data sets, however this format is not that intuitive to analysis, with this in mind we can further our work by representing this data graphically. In order to do this we choose to plot both sets of data in a boxplot and also histogram . These both give us a better 'feel' for the distribution of the data. The plots for partly vaccinated people can be found in Figure 3 and the plots for GDP per capita can be found in Figure 5.

In Figure 3, we observe that the boxplot has a large spread between the 'minimum'

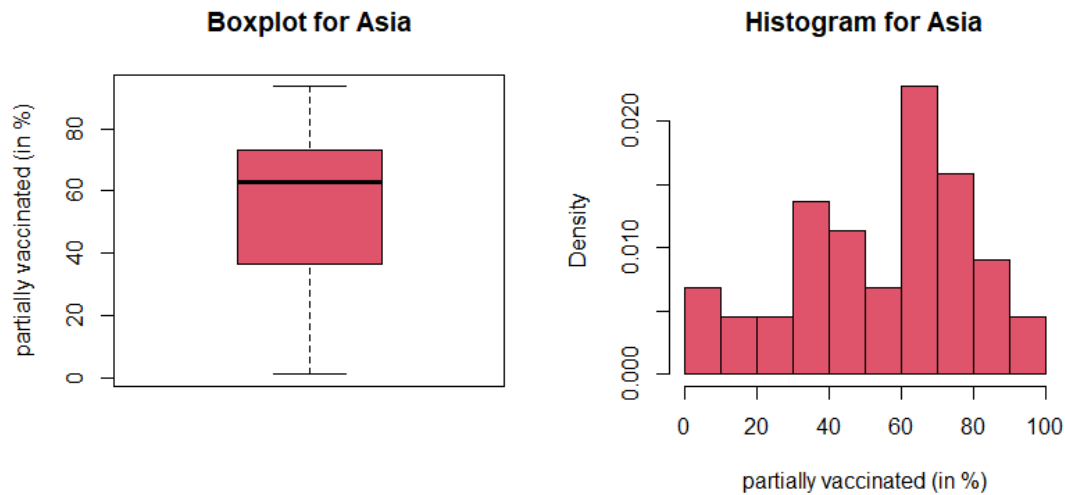


Figure 3: Graphical representation of partly vaccinated people for Asia. Left: Boxplot and Right: histogram.

and maximum', which nearly covers the whole range possible for this variable, thus there are no outliers. From the boxplot and the histogram we can see that the data is skewed to the upper side, indicated by the black median line in the boxplot and the larger bars in the histogram.

Analysing Figure 4, we see that in this case we have three outliers, namely Qatar, Macao and Singapore. We also notice that the data in this case is positively skewed, which is seen by in the boxplot with the low median and in the histogram with the higher bars closer to the origin of the x-axis.

The next step in our analysis is to consider both these univariant data sets as a bivariate data set. A very effective method of showing the relationship between two variables is graphically through a scatter plot. In Figure 5 we plot GDP against the partly vaccinated people. This figure indicates that countries with higher GDP have more people partly vaccinated. The opposite, that countries with a lower GDP have less people vaccinated is not directly true. As we can see there are some countries with a relatively low GDP having higher partly vaccinated parentage than most other countries in Asia. One may speculate this maybe due to foreign aid given to this country.

We can also do a numerical summary of the bivariate data set, this can be seen in Tables 3, 4 and 5. These are matrix which numerical represent the correlation between the GDP and the the partly vaccinated people.

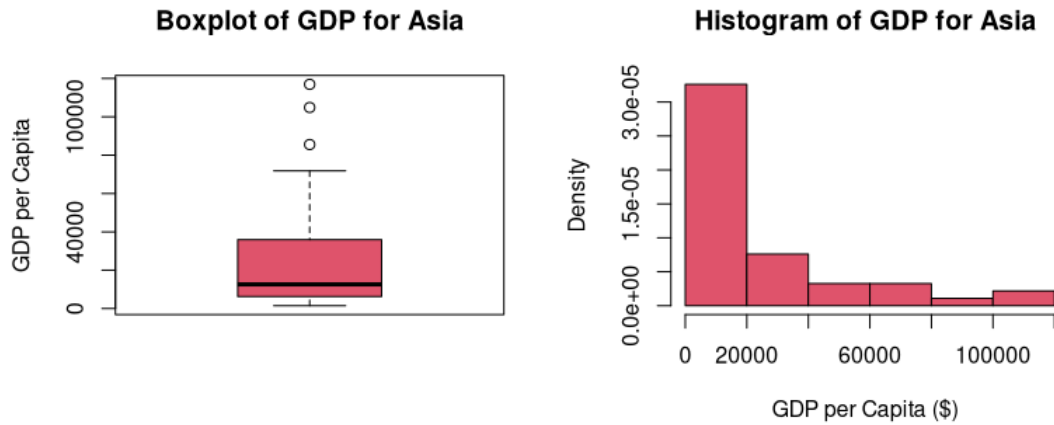


Figure 4: Graphical representation of GPA per capita for Asia. Left: boxplot and Right: histogram.

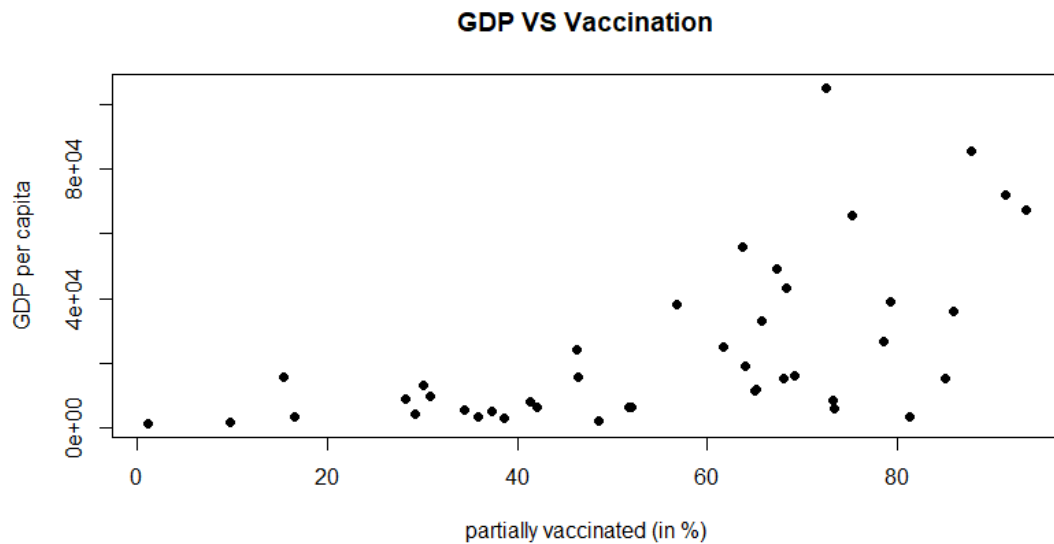


Figure 5: Scatter plot of GDP against partly vaccinated people.

Table 3: **Correlation matrix**

1.00	0.60
0.60	1.00

Table 5: **Kendall’s rank correlation coefficients**

1.00	0.47
0.47	1.00

Table 4: **Spearman’s rank correlation coefficients**

1.00	0.65
0.65	1.00

1 Appendix

R code for exercise 1.4

```
1 # Exercise 1.4
2
3 #Create function
4 lognorm <- function(n,mu,sigma){
5   # set the seed 20220211
6   # Our group number is 36
7   input_seed = 20220211 + 36
8   set.seed(input_seed)
9
10  #create vector of random
11  rnom_data = rlnorm(n, mu , sigma )
12
13  # Compute the Quantiles at %5 %10 %25 %50 %75 %90 %95.
14  quants = quantile(rnom_data, probs=c(0.05, 0.1,0.25,0.5,0.75,0.9,0.95),
15    type=7)
16
17  # Arithmetic Mean
18  loc = mean(rnom_data)
19
20  # Standard Deviation
21  spread = sd(rnom_data)
22
23  # Student numbers
24  student_nu= c(2663347, 2640428)
25
26  # put allcomputed variables intoa list
27  mylist=list( quants, loc , spread , student_nu)
28
29  # Save to disk the return of the list
30  save(mylist, file="/home/james/SDA_files/myfile1_36.RData")
}
```

code/exercise1.4.R

R code for exercise 1.5

```
1 # Exercise 1.5
2
3 bin <- function(n,size ,prob){
4   #binomial data random create
5   binomial_data = rbinom(n, size , prob)
6
7   #poisson data create
8   lambda_user= size * prob
9   poisson_data = dpois( seq(0,size ,1) , lambda_user)
10
11  #breaks
12  breaks = 0:max((7*size*prob) , max(binomial_data)+1)-0.001
13
14  #plot histogram and line function
15  { hist(binomial_data , prob=T, breaks )
16    lines(poisson_data , type="s" , col="red" ) }
17 }
```

code/exercise1.5.R

R code for exercise 1.6

```

1 #Exercise 1.6 Covid Asia Data
2
3 # Import data from file and store as variable
4 attach(covid_data)
5 covid_data <- read.csv("owid_covid_data_selection.csv", header=TRUE)
6 covid_data = owid_covid_data_selection
7
8 # Create data with only Asain
9 covid_data_asia <- covid_data[covid_data$continent == "Asia", ]
10
11
12 #####
13 #           Vaccines
14 #####
15
16 # Set up a display two pplot
17 par(mfrow=c(1,2))
18
19 # Create first plot - Box plot
20 boxplot(covid_data_asia$partly_vacc, main="Boxplot of part vac for Asia",
21         ylab="partially vaccinated (in %)", names="Asia", col=2)
22
23 # Create second plot -histgram
24 par(new=F)
25 hist(covid_data_asia$partly_vacc, main="Histogram of part vac for Asia",
26       xlab="partially vaccinated (in %)",
27       prob=T, xlim=c(0,100), col=2)
28
29
30 # Partly vaccinated asia data
31 covid_data_asia_vac = covid_data_asia[, 'partly_vacc']
32 covid_data_asia_vacNA <- covid_data_asia_vac[!is.na(covid_data_asia_vac)]
33 vac_sample_size = nrow(covid_data_asia_vac)
34 vac_NAS = sample_size - length(covid_data_asia_vacNA)
35 vac_sum = summary(covid_data_asia_vacNA)
36 vac_iqr=IQR
37
38
39 #####
40 #           GDP
41 #####
42 # Set up a display two plot
43 par(mfrow=c(1,2))
44
45 # Create first plot - Box plot
46 boxplot(covid_data_asia$gdp_per_capita, main="Boxplot of GDP for Asia",
47         ylab="GDP per Capita", names="Asia", col=2)
48
49 # Create second plot -histgram
50 par(new=F)
51 hist(covid_data_asia$gdp_per_capita, main="Histogram of GDP for Asia",
52       xlab="GDP per Capita ($)",
53       prob=T, xlim=c(0,116936), col=2)

```

```

54
55
56
57 # gdp asia data
58 covid_data_asia_gdp = covid_data_asia[, 'gdp_per_capita']
59 covid_data_asia_gdpNA <- covid_data_asia_gdp[!is.na(covid_data_asia_gdp)]
60 gdp_sample_size = nrow(covid_data_asia_gdp)
61 gdp_NAS = sample_size - length(covid_data_asia_gdp)
62 gdp_sum = summary(covid_data_asia_gdpNA)
63 gdp_iqr=IQR(covid_data_asia_gdpNA)
64 sd(covid_data_asia_gdpNA)
65
66
67
68 #####
69 #           Bivariate
70 #####
71
72 # Bivariate data for Asia: percent vaccinated and gdp
73 bivariate_asia <- cbind(covid_data_asia$partly_vacc, covid_data_asia$gdp_per_
    per_capita)
74
75 # Column means
76 colMeans(bivariate_asia, na.rm=T)
77
78 # Remove data fields tha contain NA
79 bivariate_withoutNA <- bivariate_asia[~which(is.na(covid_data_asia$gdp_per_
    capita)
80                                     | is.na(covid_data_asia$partly
    _vacc)), ]
81
82 #Correlation matrix
83 cor(bivariate_withoutNA)
84 #Spearman's rank correlation coefficients
85 cor(bivariate_withoutNA, method = "spearman")
86 # Kendall's rank correation coefficients
87 cor(bivariate_withoutNA, method = "kendall")
88
89 # Plot bivariate scatter plot
90 par(mfrow=c(1,1))
91 plot(bivariate_withoutNA[, 1], bivariate_withoutNA[, 2],
92      main="Scatterplot for Asia GDP VS Vaccination",
93      ylab="GDP per capita", xlab="partially vaccinated (in %)", pch=19)

```

code/exercise1.6.R