

Name: Ong Jing Wei

Student ID: 32909764

Generative AI was used in this assignment.

### 1. Descriptive analysis and pre-processing. (6 Marks)

(a) Describe the data overall, including things such as dimension, data types, distribution of numerical attributes, variety of non-numerical (text) attributes, missing values, and anything else of interest or relevance.

This dataset consists of **40000** rows and **52** columns.

There are 52 attributes in this dataset where **45** attributes which belongs to the **integer** class are:

"employstatus\_1", "employstatus\_2", "employstatus\_3", "employstatus\_4", "employstatus\_5", "employstatus\_6", "employstatus\_7", "employstatus\_8", "employstatus\_9", "employstatus\_10", "isoFriends\_inPerson", "isoOthPpl\_inPerson", "isoFriends\_online", "isoOthPpl\_online", "lone01", "lone02", "lone03", "happy", "lifeSat", "MLQ", "bor01", "bor02", "bor03", "consp01", "consp02", "consp03", "c19perBeh01", "c19perBeh02", "c19perBeh03", "c19RCA01", "c19RCA02", "c19RCA03", "coronaClose\_1", "coronaClose\_2", "coronaClose\_3", "coronaClose\_4", "coronaClose\_5", "coronaClose\_6", "gender", "age", "edu", "c19ProSo01", "c19ProSo02", "c19ProSo03", "c19ProSo04"

and **7** attributes which belongs to the **character** class are:

"rankOrdLife\_1", "rankOrdLife\_2", "rankOrdLife\_3", "rankOrdLife\_4", "rankOrdLife\_5", "rankOrdLife\_6", "coded\_country"

```
> summary(covid)
employstatus_1 employstatus_2 employstatus_3 employstatus_4 employstatus_5 employstatus_6 employstatus_7 employstatus_8 employstatus_9 employstatus_10
Min. :1 Min. :1
1st Qu.:1 1st Qu.:1
Median :1 Median :1
Mean :1 Mean :1 Mean :1 Mean :1 Mean :1 Mean :1 Mean :1 Mean :1 Mean :1
3rd Qu.:1 3rd Qu.:1
Max. :1 Max. :1
NA's :34362 NA's :33277 NA's :29151 NA's :36452 NA's :37972 NA's :36892 NA's :36439 NA's :39272 NA's :31782 NA's :39073
isoFriends_inPerson isoOthPpl_inPerson isoFriends_online isoOthPpl_online lone01 lone02 lone03 happy lifeSat
Min. :0.000 Min. :0.00 Min. :0.000 Min. :0.000 Min. :1.000 Min. :1.000 Min. :1.000 Min. :1.000 Min. :1.000
1st Qu.:0.000 1st Qu.:0.00 1st Qu.:2.000 1st Qu.:0.000 1st Qu.:1.000 1st Qu.:2.000 1st Qu.:1.000 1st Qu.:5.000 1st Qu.:3.000
Median :1.000 Median :1.00 Median :5.000 Median :2.000 Median :3.000 Median :2.000 Median :7.000 Median :4.000
Mean :2.073 Mean :1.96 Mean :4.397 Mean :2.866 Mean :2.418 Mean :2.664 Mean :2.078 Mean :6.332 Mean :4.136
3rd Qu.:4.000 3rd Qu.:3.00 3rd Qu.:7.000 3rd Qu.:5.000 3rd Qu.:3.000 3rd Qu.:4.000 3rd Qu.:3.000 3rd Qu.:8.000 3rd Qu.:5.000
Max. :7.000 Max. :7.00 Max. :7.000 Max. :7.000 Max. :5.000 Max. :5.000 Max. :10.000 Max. :6.000
NA's :320 NA's :520 NA's :942 NA's :1154 NA's :77 NA's :119 NA's :137 NA's :500 NA's :107
MLQ bor01 bor02 bor03 consp01 consp02 consp03 rankOrdLife_1 rankOrdLife_2
Min. :-3.000 Min. :-3.000 Min. :-3.0000 Min. :-3.0000 Min. :0.000 Min. :0.000 Min. :0.000 Length:40000 Length:40000
1st Qu.:0.000 1st Qu.:-1.000 1st Qu.:-2.0000 1st Qu.:-1.0000 1st Qu.:5.00 1st Qu.:5.000 1st Qu.:4.000 Class :character Class :character
Median :1.000 Median :0.000 Median :0.0000 Median :0.0000 Median :7.00 Median :8.000 Median :5.000 Mode :character Mode :character
Mean :0.842 Mean :0.3222 Mean :0.0421 Mean :0.3141 Mean :6.83 Mean :7.135 Mean :5.577
3rd Qu.:2.000 3rd Qu.:2.0000 3rd Qu.:2.0000 3rd Qu.:2.0000 3rd Qu.:9.00 3rd Qu.:9.000 3rd Qu.:8.000
Max. :3.000 Max. :3.0000 Max. :3.0000 Max. :10.00 Max. :10.000 Max. :10.000
NA's :109 NA's :144 NA's :164 NA's :167 NA's :1519 NA's :1544 NA's :1565

rankOrdLife_3 rankOrdLife_4 rankOrdLife_5 rankOrdLife_6 c19perBeh01 c19perBeh02 c19perBeh03 c19RCA01 c19RCA02
Length:40000 Length:40000 Length:40000 Length:40000 Min. :-3.000 Min. :-3.000 Min. :-3.000 Min. :-3.000
Class :character Class :character Class :character Class :character 1st Qu.: 2.000 1st Qu.: 2.000 1st Qu.: 1.00 1st Qu.: 0.000 1st Qu.: 2.000
Mode :character Mode :character Mode :character Mode :character Median : 3.000 Median : 3.000 Median : 2.00 Median : 3.000
Mean : 2.313 Mean : 2.424 Mean : 1.83 Mean : 1.243 Mean : 2.044
3rd Qu.: 3.000 3rd Qu.: 3.000
Max. : 3.000 Max. : 3.000 Max. : 3.000 Max. : 3.000 Max. : 3.000 Max. : 3.000 Max. : 3.000
NA's :122 NA's :126 NA's :128 NA's :121 NA's :128
c19RCA03 coronaClose_1 coronaClose_2 coronaClose_3 coronaClose_4 coronaClose_5 coronaClose_6 gender age edu
Min. :-3.000 Min. :1 Min. :1 Min. :1 Min. :1 Min. :1 Min. :1 Min. :1.00 Min. :1.000 Min. :1.000
1st Qu.: 0.000 1st Qu.:1 1st Qu.:1 1st Qu.:1 1st Qu.:1 1st Qu.:1 1st Qu.:1 1st Qu.:1.00 1st Qu.:2.000 1st Qu.:4.000
Median : 2.000 Median :1 Median :1 Median :1 Median :1 Median :1 Median :1 Median :1.00 Median :3.000 Median :5.000
Mean : 1.154 Mean :1 Mean :1 Mean :1 Mean :1 Mean :1 Mean :1.39 Mean :2.893 Mean :4.394
3rd Qu.: 3.000 3rd Qu.:1 3rd Qu.:1 3rd Qu.:1 3rd Qu.:1 3rd Qu.:1 3rd Qu.:1 3rd Qu.:2.00 3rd Qu.:4.000 3rd Qu.:5.000
Max. : 3.000 Max. :1 Max. :1 Max. :1 Max. :1 Max. :1 Max. :3.00 Max. :8.000 Max. :7.000
NA's :137 NA's :39447 NA's :38757 NA's :38443 NA's :35096 NA's :35497 NA's :10657 NA's :225 NA's :248 NA's :270
```

coded_country	c19ProSo01	c19ProSo02	c19ProSo03	c19ProSo04
Length: 40000	Min. :-3.0000	Min. :-3.0000	Min. :-3.0000	Min. :-3.000
Class : character	1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.: 0.000
Mode : character	Median : 1.0000	Median : 1.0000	Median : 1.0000	Median : 2.000
	Mean : 0.9692	Mean : 0.6672	Mean : 0.5451	Mean : 1.274
	3rd Qu.: 2.0000	3rd Qu.: 2.0000	3rd Qu.: 2.0000	3rd Qu.: 2.000
	Max. : 3.0000	Max. : 3.0000	Max. : 3.0000	Max. : 3.000
	NA's :114	NA's :129	NA's :144	NA's :142

By looking from the summary and the boxplot (Figure1.1.1, Figure1.1.2, Figure1.1.3, Figure1.1.4), we can see that the employment status attributes from **employstatus\_1** to **employstatus\_10** are binary attributes as the minimum and maximum values of 1 across these attributes suggest a binary coding scheme where 1 indicates the label of that the variable best describes the individual's employment status. The presence of missing values implies that respondents were likely instructed to select only one option among the available employment status categories, while not applicable is set to a null value. This is same with the attributes **coronaClose\_1** to **coronaClose\_6** as the minimum and maximum value of these variables are also 1, the presence of missing values implies that respondents were likely instructed to select only one option from the available categories, leaving other options unselected. Hence the missing values here are reasonable.

The attributes **isoFriends\_inPerson**, **isoOthPpl\_inPerson**, **isoFriends\_online**, **isoOthPpl\_online**, **lone01** to **lone03**, **happy**, **lifeSat**, **consp01** to **consp03** are attributes with positive range values, where higher values indicate greater agreement or preference towards the category. While the variables **MLQ**, **bor01** to **bor03**, **c19perBeh01** to **c19perBeh03**, **c19RCA01** to **c19RCA03**, **c19ProSo01** to **c19ProSo04** are attributes with the value ranging from negative to positive, where the values indicating the extent of agreement or disagreement towards the category.

**isoFriends\_inPerson**, **isoOthPpl\_inPerson**, **isoFriends\_online**, **isoOthPpl\_online** have a high number of missing values, suggesting that the respondent might not want to disclose personal information about their social interactions. Similarly, the high number of missing values in the **happy** attribute might be due to the impact of pandemic, which could make it difficult for individuals to accurately assess their happiness. Additionally, for the high number of missing values in attributes **consp01** to **consp03** may indicate concerns about potential repercussions, such as being subject to surveillance or facing criticism for expressing distrust in political leaders or institutions.

**rankOrdLife\_1** to **rankOrdLife\_6** is the categorical attributes where the values are encoded using letters A to F, each corresponding to a specific aspect of life. Higher-ranked variable indicates greater importance or significance in the respondent's lives. These variables can be used to explore correlations between these rankings and other variables.

**gender**, **age** and **edu** are also categorical attributes where it uses integer to represent distinct categories or groups. Missing values in these attributes suggest that respondents decided not to disclose their preferences or attitudes for the corresponding categories.

Var1	Freq	7	Australia	768	96	Switzerland	41	10	Bahrain	5		
105	United States of America	7007	32	Egypt	710	100	Tunisia	40	38	Georgia	5	
94	Spain	1996	62	Malaysia	547	83	Portugal	34	52	Jamaica	5	
72	Netherlands	1922	56	Kosovo	511	8	Austria	33	54	Jordan	5	
40	Greece	1837	55	Kazakhstan	487	33	El Salvador	30	106	Uruguay	5	
86	Romania	1683	77	Poland	472	66	Mexico	28	31	Ecuador	4	
46	Indonesia	1542	43	Pakistan	470	24	Colombia	26	59	Lebanon	4	
85	Republic of Serbia	1330	22	Hungary	265	70	Morocco	23	2	Albania	3	
51	Italy	1254	26	Chile	218	48	Iraq	20	18	Brunei	3	
104	United Kingdom	1214	47	Croatia	206	67	Moldova	19	25	Costa Rica	3	
101	Turkey	1107	80	Iran	206	49	Ireland	18	30	Dominican Republic	3	
37	France	1085	42	Peru	193	99	Trinidad and Tobago	18	74	Nigeria	3	
39	Germany	1052	1	Hong Kong S.A.R.	186	73	New Zealand	17	12	Belarus	2	
23	China	1006	109		157	61	Luxembourg	13	34	Estonia	2	
21	Canada	970	89	Vietnam	147	29	Denmark	12	41	Guatemala	2	
81	Philippines	963	3	Singapore	141	78	Palestine	12	57	Kuwait	2	
87	Russia	921	11	Algeria	117	28	Czech Republic	11	79	Panama	2	
88	Saudi Arabia	908	97	Bangladesh	99	36	Finland	11	84	Qatar	2	
93	South Korea	897	45	Thailand	97	63	Mali	11	107	Uzbekistan	2	
92	South Africa	896	103	Taiwan	94	15	Bosnia and Herzegovina	9	4	Andorra	1	
102	Ukraine	896	27	India	60	19	Bulgaria	9	6	44	Iceland	1
5	Argentina	868	95	United Arab Emirates	57	60	Lithuania	8	9	Armenia	1	
17	Brazil	864	50	Cyprus	51	75	Norway	8	14	58	Kyrgyzstan	1
53	Japan	835	13	Sweden	49	108	Venezuela	8	16	64	Malta	1
			Israel	48	90	Slovakia	7	20	65	Mauritius	1	
			Belgium	42	69	Montenegro	6	35	68	Mongolia	1	
									71	Myanmar	1	
									76	Oman	1	
									91	Slovenia	1	

The **coded\_country** attribute consists of a list of countries with a total of 109 distinct countries. Notably, most of the respondents are from the United States of America, with a frequency count of 7007. While Andorra, Armenia, Azerbaijan, Benin, Botswana, Cameroon, Ethiopia, Iceland, Kyrgyzstan, Malta, Mauritius, Mongolia, Myanmar, Oman, and Slovenia have the least frequency which is only 1. However, there is an empty value observed between Hong Kong S.A.R and Vietnam, occurring 157 times. This suggests that there might be some incomplete values for this attribute, or the respondent might not want to disclose where they are from.

(b) Comment on any pre-processing or data manipulation required for the following analysis.

For data pre-processing, I addressed missing values in several attributes. Initially, I replaced NULL values in **employstatus\_1** to **employstatus\_10** and **coronaClose\_1** to **coronaClose\_6** with 0. This adjustment was necessary as respondents were instructed to select only one option from available categories, leaving others unselected. To avoid dropping entire rows due to these attributes, I replaced NA values accordingly. Additionally, I converted empty values in the **coded\_country** attribute to NA. Since empty values do not contribute to our analysis, replacing them with NA facilitated their removal in subsequent steps. After handling missing values, I performed data clean up by dropping rows with NA values, preparing the dataset for further analysis.

Furthermore, I encoded categorical attributes **rankOrdLife\_1** to **rankOrdLife\_6** into numeric values for easier comparison. This involved defining a function to encode categorical variables and applying it to each column using a mapping scheme.

## 2. Focus country vs all other countries as a group. (12 Marks)

(a) Identify your focus country from the accompanying list (FocusCountryByID.pdf). How do participant responses (attributes) for your focus country differ from the other countries in the survey as a group?

When comparing the summary statistics below and the boxplot (Figure 2.1.1, Figure 2.1.2, Figure 2.1.3) between the Netherlands and other countries, we observe distinct patterns. In the Netherlands, a significant proportion of respondents are employed, and working 24-39 hours per week with around 75% as indicated by the 3rd quartile of **employstatus\_2**. The mean of 0.3001 suggests that approximately 30% of respondents fall into this category, reflecting a common employment pattern. Conversely, in other countries, respondents are predominantly working 40 or more hours per week.

Moreover, respondents from the Netherlands appear to have more social interactions with individuals outside their households compared to those from other countries. This is evident from higher mean, median, and 3rd quartile values in **isoOthPpl\_inPerson**. Additionally, the p-value for the t-test below is extremely small also indicates that there is significantly larger mean value of **isoOthPpl\_inPerson** in the Netherlands compared to other countries, further highlighting this distinction.

```

> # t-test for iso0thPpl_inPerson
> t_test_iso0thPpl_inPerson <- t.test(netherlands_data$iso0thPpl_inPerson, other_countries_data$iso0thPpl_inPerson)
> t_test_iso0thPpl_inPerson

Welch Two Sample t-test

data: netherlands_data$iso0thPpl_inPerson and other_countries_data$iso0thPpl_inPerson
t = 10.883, df = 1538.9, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
0.5575795 0.8027546
sample estimates:
mean of x mean of y
2.525471 1.845304

```

Respondents from Netherlands show lower tendencies towards feelings of isolation and exclusion, as indicated by lower median and mean values in lone02 and lone03 compared to respondents from other countries. Similarly, the mean values for conspiracy beliefs (consp01, consp02, consp03) in the Netherlands are lower than in other countries, suggesting a lower inclination towards believing in conspiracies. This implies that individuals in the Netherlands may have higher levels of trust in public information, political transparency, and government agencies' actions compared to respondents from other countries. A t-test conducted for one of the conspiracy attributes (consp02) reveals an extremely small p-value, indicating a significant difference in mean values between the Netherlands and other countries.

```

> # t-test for consp02
> t_test_consp02 <- t.test(netherlands_data$consp02, other_countries_data$consp02)
> t_test_consp02

Welch Two Sample t-test

data: netherlands_data$consp02 and other_countries_data$consp02
t = -24.907, df = 1546, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-1.846184 -1.576630
sample estimates:
mean of x mean of y
5.625262 7.336668

```

```

> summary(netherlands_data)

employstatus_1 employstatus_2 employstatus_3 employstatus_4 employstatus_5 employstatus_6
Min. :0.0000 Min. :0.0000 Min. :0.0000 Min. :0.00000 Min. :0.00000 Min. :0.00000
1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.00000 1st Qu.:0.00000 1st Qu.:0.00000
Median :0.0000 Median :0.0000 Median :0.0000 Median :0.00000 Median :0.00000 Median :0.00000
Mean :0.1445 Mean :0.3001 Mean :0.1961 Mean :0.05024 Mean :0.02442 Mean :0.06909
3rd Qu.:0.0000 3rd Qu.:1.0000 3rd Qu.:0.0000 3rd Qu.:0.00000 3rd Qu.:0.00000 3rd Qu.:0.00000
Max. :1.0000 Max. :1.0000 Max. :1.0000 Max. :1.00000 Max. :1.00000 Max. :1.00000

employstatus_7 employstatus_8 employstatus_9 employstatus_10 isoFriends_inPerson iso0thPpl_inPerson
Min. :0.0000 Min. :0.0000 Min. :0.0000 Min. :0.00000 Min. :0.00000 Min. :0.00000
1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.00000 1st Qu.:0.00000 1st Qu.:0.00000
Median :0.0000 Median :0.0000 Median :0.0000 Median :0.00000 Median :1.0000 Median :2.0000
Mean :0.1117 Mean :0.05652 Mean :0.1689 Mean :0.04676 Mean :1.909 Mean :2.525
3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:0.00000 3rd Qu.:3.000 3rd Qu.:4.000
Max. :1.0000 Max. :1.00000 Max. :1.0000 Max. :1.00000 Max. :7.000 Max. :7.000

isoFriends_online iso0thPpl_online lone01 lone02 lone03 happy
Min. :0.00 Min. :0.0000 Min. :1.000 Min. :1.00000 Min. :1.00000 Min. :1.00000
1st Qu.:2.00 1st Qu.:0.0000 1st Qu.:1.000 1st Qu.:1.00000 1st Qu.:1.00000 1st Qu.:6.000
Median :4.00 Median :2.0000 Median :2.0000 Median :2.00000 Median :1.0000 Median :7.000
Mean :3.88 Mean :2.522 Mean :1.999 Mean :2.397 Mean :1.78 Mean :6.871
3rd Qu.:7.00 3rd Qu.:5.0000 3rd Qu.:3.000 3rd Qu.:3.00000 3rd Qu.:2.000 3rd Qu.:8.000
Max. :7.00 Max. :7.0000 Max. :5.000 Max. :5.00000 Max. :5.00 Max. :10.000

lifeSat MLQ bor01 bor02 bor03 consp01
Min. :1.000 Min. :-3.0000 Min. :-3.0000 Min. :-3.00000 Min. :-3.00000 Min. :0.000
1st Qu.:4.000 1st Qu.:0.0000 1st Qu.:-2.0000 1st Qu.:-2.00000 1st Qu.:-1.00000 1st Qu.:4.000
Median :5.000 Median :1.0000 Median :0.0000 Median :0.00000 Median :1.00000 Median :6.000
Mean :4.574 Mean :0.8255 Mean :-0.0321 Mean :-0.3043 Mean :0.4955 Mean :5.846
3rd Qu.:5.000 3rd Qu.:2.0000 3rd Qu.:2.0000 3rd Qu.:1.00000 3rd Qu.:2.00000 3rd Qu.:8.000
Max. :6.000 Max. :3.0000 Max. :3.0000 Max. :3.00000 Max. :3.00000 Max. :10.000

consp02 consp03 rankOrdLife_1 rankOrdLife_2 rankOrdLife_3 rankOrdLife_4
Min. :0.000 Min. :0.0000 Min. :1.000 Min. :1.00000 Min. :1.00000 Min. :1.00000
1st Qu.:4.000 1st Qu.:2.0000 1st Qu.:4.0000 1st Qu.:3.00000 1st Qu.:4.00000 1st Qu.:2.000
Median :6.000 Median :5.0000 Median :5.0000 Median :4.00000 Median :5.00000 Median :2.000
Mean :5.625 Mean :4.418 Mean :4.385 Mean :3.999 Mean :4.861 Mean :2.436
3rd Qu.:8.000 3rd Qu.:6.0000 3rd Qu.:6.0000 3rd Qu.:5.00000 3rd Qu.:6.00000 3rd Qu.:3.000
Max. :10.000 Max. :10.0000 Max. :6.0000 Max. :6.00000 Max. :6.00000 Max. :6.000

rankOrdLife_5 rankOrdLife_6 c19perBeh01 c19perBeh02 c19perBeh03 c19RCA01
Min. :1.000 Min. :1.0000 Min. :-3.000 Min. :-3.00000 Min. :-3.00000 Min. :-3.000
1st Qu.:1.000 1st Qu.:2.0000 1st Qu.:2.0000 1st Qu.:2.00000 1st Qu.:1.00000 1st Qu.:0.000
Median :2.000 Median :3.0000 Median :2.0000 Median :3.00000 Median :2.00000 Median :2.000
Mean :2.149 Mean :3.171 Mean :2.112 Mean :2.373 Mean :1.278 Mean :0.993
3rd Qu.:3.000 3rd Qu.:4.0000 3rd Qu.:3.0000 3rd Qu.:3.00000 3rd Qu.:2.00000 3rd Qu.:3.000
Max. :6.000 Max. :6.0000 Max. :3.0000 Max. :3.00000 Max. :3.00000 Max. :3.000

c19RCA02 c19RCA03 coronaClose_1 coronaClose_2 coronaClose_3 coronaClose_4
Min. :-3.000 Min. :-3.0000 Min. :0.0000 Min. :0.00000 Min. :0.00000 Min. :0.000
1st Qu.:1.000 1st Qu.:-2.0000 1st Qu.:0.0000 1st Qu.:0.00000 1st Qu.:0.00000 1st Qu.:0.000
Median :2.000 Median :0.0000 Median :0.00000 Median :0.00000 Median :0.00000 Median :0.000
Mean :1.683 Mean :0.1047 Mean :0.0307 Mean :0.0642 Mean :0.0635 Mean :0.164
3rd Qu.:3.000 3rd Qu.:2.0000 3rd Qu.:0.0000 3rd Qu.:0.00000 3rd Qu.:0.00000 3rd Qu.:0.000
Max. :3.000 Max. :3.0000 Max. :1.0000 Max. :1.00000 Max. :1.0000 Max. :1.000

coronaClose_5 coronaClose_6 gender age edu coded_country
Min. :0.0000 Min. :0.0000 Min. :1.000 Min. :1.00000 Min. :1.00 Length:1433
1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:1.000 1st Qu.:2.0000 1st Qu.:3.00 Class :character
Median :0.0000 Median :1.0000 Median :1.000 Median :3.00000 Median :4.00 Mode :character
Mean :0.1486 Mean :0.6315 Mean :1.369 Mean :3.338 Mean :4.39
3rd Qu.:0.0000 3rd Qu.:1.0000 3rd Qu.:2.000 3rd Qu.:5.0000 3rd Qu.:6.00
Max. :1.0000 Max. :1.0000 Max. :3.000 Max. :8.000 Max. :7.00

c19ProSo01 c19ProSo02 c19ProSo03 c19ProSo04
Min. :-3.0000 Min. :-3.0000 Min. :-3.00000 Min. :-3.000
1st Qu.:0.0000 1st Qu.:-1.0000 1st Qu.:0.0000 1st Qu.:0.000
Median :1.0000 Median :0.0000 Median :1.00000 Median :2.000
Mean :0.9525 Mean :0.2114 Mean :0.6364 Mean :1.337
3rd Qu.:2.0000 3rd Qu.:2.0000 3rd Qu.:2.0000 3rd Qu.:2.000
Max. :3.0000 Max. :3.0000 Max. :3.00000 Max. :3.000

```

```

> summary(other_countries_data)
employstatus_1 employstatus_2 employstatus_3 employstatus_4 employstatus_5 employstatus_6
Min. :0.0000 Min. :0.0000 Min. :0.0000 Min. :0.0000 Min. :0.0000 Min. :0.00000
1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.00000
Median :0.0000 Median :0.0000 Median :0.0000 Median :0.0000 Median :0.0000 Median :0.00000
Mean :0.1424 Mean :0.1608 Mean :0.2767 Mean :0.09127 Mean :0.0528 Mean :0.07864
3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:1.0000 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:0.00000
Max. :1.0000 Max. :1.0000 Max. :1.0000 Max. :1.0000 Max. :1.0000 Max. :1.00000
employstatus_7 employstatus_8 employstatus_9 employstatus_10 isoFriends_inPerson
Min. :0.000 Min. :0.0000 Min. :0.0000 Min. :0.00000 Min. :0.000
1st Qu.:0.000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.00000 1st Qu.:0.000
Median :0.000 Median :0.0000 Median :0.0000 Median :0.00000 Median :1.000
Mean :0.091 Mean :0.0162 Mean :0.2103 Mean :0.02272 Mean :1.966
3rd Qu.:0.000 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:0.00000 3rd Qu.:3.000
Max. :1.000 Max. :1.0000 Max. :1.0000 Max. :1.00000 Max. :7.000
iso0thPpl_inPerson isoFriends_online iso0thPpl_online lone01 lone02 lone03
Min. :0.000 Min. :0.000 Min. :0.000 Min. :1.000 Min. :1.0 Min. :1.00
1st Qu.:0.000 1st Qu.:2.000 1st Qu.:0.0000 1st Qu.:1.000 1st Qu.:2.0 1st Qu.:1.00
Median :1.000 Median :5.000 Median :2.000 Median :2.000 Median :3.0 Median :2.00
Mean :1.845 Mean :4.454 Mean :2.867 Mean :2.436 Mean :2.7 Mean :2.08
3rd Qu.:3.000 3rd Qu.:7.000 3rd Qu.:5.000 3rd Qu.:3.000 3rd Qu.:4.0 3rd Qu.:3.00
Max. :7.000 Max. :7.000 Max. :7.000 Max. :5.000 Max. :5.0 Max. :5.00
    happy lifeSat MLQ bor01 bor02
Min. : 1.000 Min. :1.000 Min. :-3.000 Min. :-3.0000 Min. :-3.0000
1st Qu.: 5.000 1st Qu.:3.000 1st Qu.: 0.000 1st Qu.:-1.0000 1st Qu.:-2.00000
Median : 7.000 Median :4.000 Median : 1.000 Median : 0.0000 Median : 0.00000
Mean : 6.331 Mean :4.145 Mean : 0.867 Mean : 0.3217 Mean : 0.05991
3rd Qu.: 8.000 3rd Qu.:5.000 3rd Qu.: 2.000 3rd Qu.: 2.0000 3rd Qu.: 2.00000
Max. :10.000 Max. :6.000 Max. : 3.000 Max. : 3.0000 Max. : 3.00000
    bor03 consp01 consp02 consp03 rankOrdLife_1 rankOrdLife_2
Min. : -3.0000 Min. : 0.000 Min. : 0.000 Min. : 0.000 Min. :1.000 Min. :1.000
1st Qu.:-1.0000 1st Qu.: 5.000 1st Qu.: 6.000 1st Qu.: 4.000 1st Qu.:4.000 1st Qu.:3.000
Median : 0.0000 Median : 7.000 Median : 8.000 Median : 6.000 Median :5.000 Median :4.000
Mean : 0.2955 Mean : 6.992 Mean : 7.337 Mean : 5.691 Mean : 4.659 Mean :3.531
3rd Qu.: 2.0000 3rd Qu.: 9.000 3rd Qu.:10.000 3rd Qu.: 8.000 3rd Qu.:6.000 3rd Qu.:5.000
Max. : 3.0000 Max. :10.000 Max. :10.000 Max. :10.000 Max. :6.000 Max. :6.000

rankOrdLife_3 rankOrdLife_4 rankOrdLife_5 rankOrdLife_6 c19perBeh01 c19perBeh02
Min. :1.000 Min. :1.000 Min. :1.000 Min. :1.000 Min. :-3.000 Min. :-3.000
1st Qu.:4.000 1st Qu.:2.000 1st Qu.:1.000 1st Qu.:2.000 1st Qu.: 2.000 1st Qu.: 2.000
Median :5.000 Median :3.000 Median :2.000 Median :3.000 Median : 3.000 Median : 3.000
Mean :4.748 Mean :2.725 Mean :2.098 Mean :3.239 Mean : 2.372 Mean : 2.483
3rd Qu.:6.000 3rd Qu.:3.000 3rd Qu.:3.000 3rd Qu.:4.000 3rd Qu.: 3.000 3rd Qu.: 3.000
Max. :6.000 Max. :6.000 Max. :6.000 Max. :6.000 Max. : 3.000 Max. : 3.000
    c19perBeh03 c19RCA01 c19RCA02 c19RCA03 coronaClose_1 coronaClose_2
Min. : -3.000 Min. : -3.000 Min. : -3.000 Min. : -3.000 Min. :0.00000 Min. :0.00000
1st Qu.: 1.000 1st Qu.: 0.000 1st Qu.: 2.000 1st Qu.: 0.000 1st Qu.:0.00000 1st Qu.:0.00000
Median : 2.000 Median : 2.000 Median : 3.000 Median : 2.000 Median :0.00000 Median :0.00000
Mean : 1.894 Mean : 1.296 Mean : 2.122 Mean : 1.188 Mean : 0.01033 Mean : 0.02862
3rd Qu.: 3.000 3rd Qu.: 3.000 3rd Qu.: 3.000 3rd Qu.: 3.000 3rd Qu.:0.00000 3rd Qu.:0.00000
Max. : 3.000 Max. : 3.000 Max. : 3.000 Max. : 3.000 Max. : 1.00000 Max. : 1.00000
coronaClose_3 coronaClose_4 coronaClose_5 coronaClose_6 gender age
Min. :0.00000 Min. :0.00000 Min. :0.00000 Min. :0.00000 Min. :1.000 Min. :1.000
1st Qu.:0.00000 1st Qu.:0.00000 1st Qu.:0.00000 1st Qu.:0.00000 1st Qu.:1.000 1st Qu.:2.000
Median :0.00000 Median :0.00000 Median :0.00000 Median :1.00000 Median :1.000 Median :3.000
Mean :0.03642 Mean :0.1236 Mean :0.1119 Mean :0.7411 Mean :1.386 Mean :2.901
3rd Qu.:0.00000 3rd Qu.:0.00000 3rd Qu.:0.00000 3rd Qu.:1.00000 3rd Qu.:2.000 3rd Qu.:4.000
Max. :1.00000 Max. :1.00000 Max. :1.00000 Max. :1.00000 Max. :3.000 Max. :8.000
    edu coded_country c19ProSo01 c19ProSo02 c19ProSo03
Min. :1.000 Length:33582 Min. : -3.000 Min. : -3.000 Min. : -3.000
1st Qu.:4.000 Class :character 1st Qu.: 0.0000 1st Qu.: 0.0000 1st Qu.: 0.000
Median :5.000 Mode :character Median : 1.0000 Median : 1.0000 Median : 1.000
Mean :4.404 Mean : 0.9845 Mean : 0.6743 Mean : 0.541
3rd Qu.:5.000 3rd Qu.: 2.0000 3rd Qu.: 2.0000 3rd Qu.: 2.000
Max. :7.000 Max. : 3.0000 Max. : 3.0000 Max. : 3.000
    c19ProSo04
Min. : -3.000
1st Qu.: 1.000
Median : 2.000
Mean : 1.339
3rd Qu.: 3.000
Max. : 3.000

```

(b) How well do participant responses (attributes) predict pro-social attitudes (**c19ProSo01, 2, 3 and 4**) for your focus country? Which attributes seem to be the best predictors? Explain your reasoning.

<b>c19ProSo01</b>	
1 <sup>st</sup> Best Predictor	c19perBeh02, p-value: 0.0000621
2 <sup>nd</sup> Best Predictor	isoOthPpl_inPerson, p-value: 0.000678
R-squared: 0.3381	RSE: 1.2
<b>c19ProSo02</b>	
1 <sup>st</sup> Best Predictor	c19RCA01, p-value: 0.000000579
2 <sup>nd</sup> Best Predictor	employstatus_7, p-value: 0.0000903
R-squared: 0.2912	RSE: 1.443
<b>c19ProSo03</b>	
1 <sup>st</sup> Best Predictor	age, p-value: 0.000000344
2 <sup>nd</sup> Best Predictor	Employstatus_9, p-value: 0.02640
R-squared: 0.4557	RSE: 1.306
<b>c19ProSo04</b>	
1 <sup>st</sup> Best Predictor	c19perBeh03, p-value: 0.000000418
2 <sup>nd</sup> Best Predictor	employstatus_9, p-value: 0.000342
R-squared: 0.442	RSE: 1.17

The table above summarizes the best predictors for all pro-social attitudes based on the linear regression analysis. These predictors were chosen because of their significantly low p-values, suggesting a strong correlation with the response variable. This indicates compelling evidence of their influence on the respondents' pro-social attitudes.

The average R-squared value, around 0.38, suggests that the model accounts for a moderate portion of the variability in the response variable. While this indicates some level of usefulness, it also implies that the model may not fully capture all underlying relationships in the data.

Furthermore, the low average residual standard error, approximately 1.28, implies that the model fits the data well and can make accurate predictions. This suggests that the model's predictions are generally close to the actual values, indicating its reliability in predicting pro-social attitudes based on the chosen predictors.

Call: lm(formula = c19ProSo01 ~ ., data = netherlands_data)	Call: lm(formula = c19ProSo02 ~ ., data = netherlands_data)
Residuals:	Residuals:
Min 1Q Median 3Q Max -5.1439 -0.6452 0.1016 0.7854 3.5074	Min 1Q Median 3Q Max -4.6165 -0.9706 0.1850 1.0320 3.9915
Coefficients: (1 not defined because of singularities)	Coefficients: (1 not defined because of singularities)
(Intercept) -0.3739872 0.5480921 -0.682 0.495136 employstatus_1 -0.1412853 0.1201039 -1.176 0.239654 employstatus_2 -0.0258750 0.1272813 -0.203 0.838938 employstatus_3 -0.0244965 0.1406741 -0.174 0.861783 employstatus_4 -0.0304483 0.1677422 -0.182 0.855987 employstatus_5 0.0927445 0.2189730 0.424 0.671965 employstatus_6 -0.1447712 0.1425236 -1.016 0.309916 employstatus_7 -0.3633045 0.1623902 -2.237 0.025430 * employstatus_8 -0.3035025 0.1720180 -1.764 0.077891 . employstatus_9 0.2416462 0.1335435 1.809 0.070591 . employstatus_10 0.1607027 0.1572076 1.022 0.306850 isoFriends_inPerson 0.0218011 0.0169090 1.289 0.197502 iso0thPpl_inPerson 0.0543594 0.0159595 3.406 0.000678 *** isoFriends_online 0.0128673 0.0147736 0.871 0.383924 iso0thPpl_online -0.0167216 0.0143444 -1.166 0.243929 lone01 0.0924839 0.0478659 1.932 0.053545 . lone02 0.0147952 0.0406912 0.364 0.716215 lone03 -0.0656983 0.0440279 -1.492 0.135875 happy -0.0400213 0.0318541 -1.256 0.209186 lifeSat 0.0304800 0.0505524 0.603 0.546648 MLQ -0.0106343 0.0269340 -0.395 0.693031 bor01 0.0020169 0.0250294 0.081 0.935786 bor02 -0.0094411 0.0259699 -0.364 0.716257 bor03 0.0134239 0.0214530 0.626 0.531591 consp01 0.0142699 0.0176812 0.807 0.419767 consp02 0.0233618 0.0191450 1.220 0.222576 consp03 -0.0041894 0.0156625 -0.267 0.789141 rankOrdLife_1 -0.0276882 0.0296229 -0.935 0.350112 rankOrdLife_2 0.0006791 0.0326239 0.021 0.983397 rankOrdLife_3 -0.0091495 0.0336959 -0.272 0.786023 rankOrdLife_4 -0.0094076 0.0378821 -0.248 0.803910 rankOrdLife_5 -0.0298346 0.0329224 -0.906 0.364984 rankOrdLife_6 NA NA NA NA c19perBeh01 0.0392345 0.0343554 1.142 0.253644	(Intercept) -3.691797 0.651449 -5.667 1.77e-08 *** employstatus_1 0.293672 0.144232 2.036 0.041930 * employstatus_2 0.375359 0.152673 2.459 0.014071 * employstatus_3 0.500717 0.168568 2.970 0.003025 ** employstatus_4 0.169070 0.201593 0.839 0.401800 employstatus_5 0.022567 0.263242 0.086 0.931696 employstatus_6 0.167909 0.171331 0.980 0.327245 employstatus_7 0.763722 0.194480 3.927 9.03e-05 *** employstatus_8 0.040813 0.207011 0.197 0.843737 employstatus_9 -0.017983 0.160721 -0.112 0.910926 employstatus_10 -0.122325 0.189021 -0.647 0.517641 isoFriends_inPerson 0.028805 0.020324 1.417 0.156612 iso0thPpl_inPerson -0.019017 0.019258 -0.987 0.323590 isoFriends_online 0.016121 0.017759 0.908 0.364156 iso0thPpl_online 0.038127 0.017221 2.214 0.026994 * lone01 0.098043 0.057557 1.703 0.088713 . lone02 -0.092609 0.048854 -1.896 0.058216 . lone03 0.102982 0.052896 1.947 0.051750 . happy 0.039018 0.038299 1.019 0.308494 lifeSat 0.109931 0.060705 1.811 0.070372 . MLQ 0.064929 0.032332 2.008 0.044817 * bor01 0.094382 0.029980 3.148 0.001678 ** bor02 0.003344 0.031220 0.107 0.914710 bor03 0.018609 0.025787 0.722 0.470649 consp01 -0.049076 0.021219 -2.313 0.020874 * consp02 -0.035843 0.023006 -1.558 0.119467 consp03 0.012170 0.018825 0.646 0.518095 rankOrdLife_1 0.090538 0.035537 2.548 0.010951 * rankOrdLife_2 0.134876 0.039049 3.454 0.000569 *** rankOrdLife_3 0.102948 0.040412 2.547 0.010958 * rankOrdLife_4 0.067423 0.045503 1.482 0.138640 rankOrdLife_5 0.122247 0.039451 3.099 0.001983 ** rankOrdLife_6 NA NA NA NA c19perBeh01 0.039476 0.041304 0.956 0.339372
c19perBeh02 0.1793836 0.0446541 4.017 6.21e-05 *** c19perBeh03 -0.0742231 0.0251008 -2.957 0.003159 ** c19RCA01 0.0165960 0.0201593 0.823 0.410512 c19RCA02 0.0055002 0.0292648 0.188 0.850948 c19RCA03 -0.0399355 0.0209970 -1.902 0.057383 . coronaClose_1 0.3801891 0.2226826 1.707 0.087988 . coronaClose_2 -0.0882258 0.1616685 -0.546 0.585346 coronaClose_3 -0.0342230 0.1539486 -0.222 0.824112 coronaClose_4 0.2792447 0.1311729 2.129 0.033445 * coronaClose_5 0.1174055 0.1357323 0.865 0.387201 coronaClose_6 0.0254124 0.1424683 0.178 0.858456 gender 0.0991188 0.0685773 1.445 0.148584 age 0.0562947 0.0321358 1.752 0.080034 . edu 0.0270482 0.0253371 1.068 0.285918 c19ProSo02 0.1261517 0.0221105 5.706 1.42e-08 *** c19ProSo03 0.3341830 0.0230192 14.518 < 2e-16 *** c19ProSo04 0.0575319 0.0275390 2.089 0.036881 *	c19perBeh02 -0.013470 0.053989 -0.249 0.803018 c19perBeh03 0.046650 0.030243 1.543 0.123178 c19RCA01 0.120630 0.024021 5.022 5.79e-07 *** c19RCA02 0.026382 0.035172 0.750 0.453330 c19RCA03 0.032869 0.025258 1.301 0.193365 coronaClose_1 -0.179933 0.267924 -0.672 0.501962 coronaClose_2 0.121454 0.194334 0.625 0.532089 coronaClose_3 0.234169 0.184957 1.266 0.205699 coronaClose_4 -0.413784 0.157548 -2.626 0.008725 ** coronaClose_5 -0.392819 0.162865 -2.412 0.015998 * coronaClose_6 -0.303484 0.171068 -1.774 0.076274 . gender -0.053925 0.082486 -0.654 0.513381 age 0.055855 0.038644 1.445 0.148578 edu 0.115379 0.030312 3.806 0.000147 *** c19ProSo01 0.182293 0.031950 5.706 1.42e-08 *** c19ProSo03 0.086675 0.029613 2.927 0.003480 ** c19ProSo04 0.107079 0.033031 3.242 0.001216 **
---	---
Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1	Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1
Residual standard error: 1.2 on 1383 degrees of freedom Multiple R-squared: 0.3381, Adjusted R-squared: 0.3147 F-statistic: 14.42 on 49 and 1383 DF, p-value: < 2.2e-16	Residual standard error: 1.443 on 1383 degrees of freedom Multiple R-squared: 0.2912, Adjusted R-squared: 0.2661 F-statistic: 11.6 on 49 and 1383 DF, p-value: < 2.2e-16



(c) Repeat Question 2(b) for the other countries as a group. Which attributes are the strongest predictors? How do these attributes compare to those of your focus country?

<b>c19ProSo01</b>		
Best Predictors	bor03, gender, employstatus_10, isoOthPpl_inPerson, lone01, lone02, employstatus_4, employstatus_7, isoFriends_online, MLQ, c19perBeh01, c19perBeh02, c19perBeh03, age	
R-squared: 0.3535	RSE: 1.173	
<b>c19ProSo02</b>		
Best Predictors	employstatus_4, employstatus_5, employstatus_8, isoFriends_inPerson, isoFriends_online, lone01, lone02, lifeSat, MLQ, bor01, consp01, consp02, c19perBeh01, c19perBeh03, c19RCA01, c19RCA03, coronaClose_5, gender, edu	
R-squared: 0.394	RSE: 1.279	
<b>c19ProSo03</b>		
Best Predictors	employstatus_7, isoFriends_online, isoOthPpl_online, lone03, lifeSat, consp02, c19RCA03, coronaClose_1, age, edu	
R-squared: 0.454	RSE: 1.232	
<b>c19ProSo04</b>		
Best Predictors	employstatus_3, employstatus_8, isoFriends_inPerson, isoFriends_online, isoOthPpl_online, lone01, lone02, lifeSat, bor01, bor02, bor03, consp01, c19perBeh01, c19perBeh02, c19perBeh03, c19RCA02, c19RCA03, coronaClose_3, age	
R-squared: 0.3237	RSE: 1.274	

The table above summarizes the best predictors for all pro-social attitudes in other countries based on the linear regression analysis. Notably, all the best predictors identified for our focus country are also found to be the best predictors for other countries. Additionally, a greater number of predictors with extremely low p-values are observed in other countries compared to our focus country. The predictors for other countries capture a wider array of factors, suggesting a more diverse set of influences on pro-social attitudes that may vary across different cultural and societal contexts. The overall 4 models have a significantly small p-values which is <2.2e-16, indicating strong evidence against the null hypothesis and strong associations between predictors and the response variable.

The average R-squared value for other countries is similar to that of our focus country, Netherlands, suggesting a similar level of explained variability in the response variable. However, it's important to note that this level of explanation may not capture all the complexities of the underlying relationships in the data.

Moreover, the average residual standard error (RSE) for other countries is slightly smaller than that of our focus country, indicating a slightly better fit of the model to the data. This suggests that the model's predictions for pro-social attitudes in other countries are generally more accurate.





### 3. Focus country vs cluster of similar countries. (10 Marks)

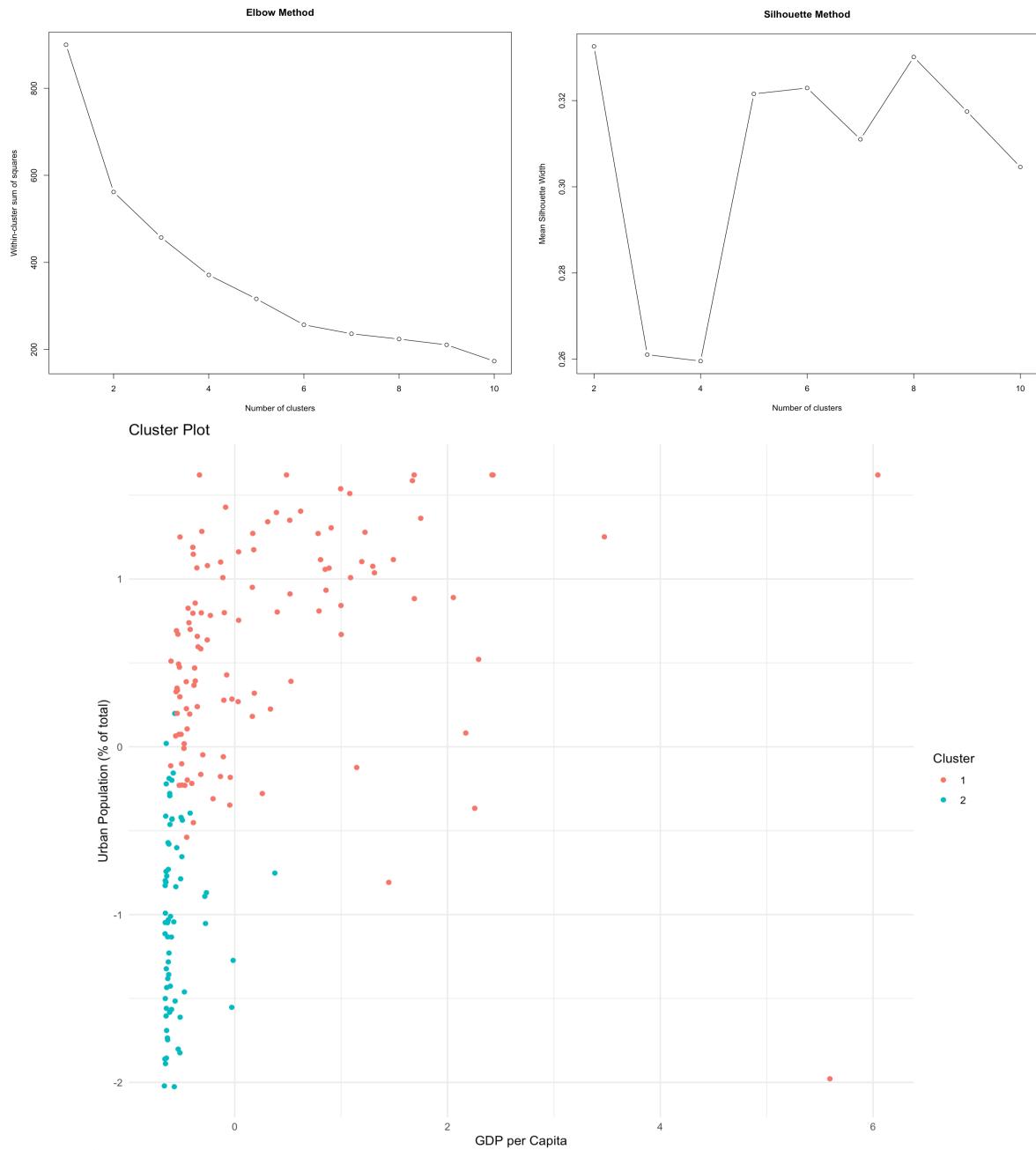
(a) Using a collection of social, economic, health, political or other indicators from external sources, identify at least 5 countries in the baseline data that are similar to your focus country using clustering. Van Lissa (2022) refers to several indicators you might consider, among others. Some of these are listed in the references, but these are not exhaustive. State the indicators used and describe how you calculated/identified similar countries. Copy and paste the table of values you used for your clustering into your report as an Appendix.

Indicator	Variable Name	Description
GDP per capita (current US\$)	latest.value_NY.GDP.PCAP.CD	This indicator measures the average economic output per person in a country, calculated by dividing the country's gross domestic product (GDP) by its population.
Urban Population (% of total population)	latest.value_SP.URB.TOTL.IN.ZS	This indicator represents the percentage of a country's total population that resides in urban areas.
Rural population (% of total population)	latest.value_SP.RUR.TOTL.ZS	This indicator represents the percentage of a country's total population that resides in rural areas.
Birth rate	latest.value_birthrate.2054	The birth rate refers to the number of live births per 1,000 people in a given population within a specified period, usually one year.
Death rate	latest.value_deathrate.2066	The death rate, also known as the mortality rate, refers to the number of deaths per 1,000 people in a given population within a specified period, usually one year.

These indicators have been chosen for their relevance in assessing the economic and demographic characteristics of countries, aiding in the identification of similar nations to our focus country. GDP per capita serves as a key measure of economic prosperity and living standards, while the percentage of urban and rural populations sheds light on a country's urbanization and agricultural landscape, respectively. Birth and death rates provide insights into demographic trends, fertility patterns, and overall population health.

To determine the optimal number of clusters for our analysis, I utilized the elbow and silhouette methods. These techniques pointed towards 2 as the ideal number of clusters. Subsequently, I conducted k-means clustering with this optimal value and identified countries sharing similarities with our focus country, Netherlands. I calculated distances between each country which belongs to the same cluster as Netherlands using Euclidean distance and sorted the data frame based on distances and get the closest countries with Netherlands.

The 6 countries identified as similar to Netherlands include **Sweden, Greenland, Austria, Australia, Finland** and **San Marino**, all falling within the same cluster. These countries exhibit comparable economic and demographic characteristics, forming a cluster that shares resemblances with Netherlands in terms of the selected indicators.



(b) How well do participant responses (attributes) predict pro-social attitudes (**c19ProSo01, 2, 3 and 4**) for this cluster of similar countries? Which attributes are the strongest predictors? How do these attributes compare to those of your focus country? Comment on the similarity and/or difference between your results for this question and Question 2(c). That is, does the group of all other countries 2(c), or the cluster of similar countries 3(b) give a better match to the important attributes for predicting pro-social attitudes in your focus country? Discuss.

<b>c19ProSo01</b>	
Best Predictors	c19perBeh01, p-value: 0.00116
	bor03, p-value: 0. 00117
R-squared: 0.4398	RSE: 1.134
<b>c19ProSo02</b>	
Best Predictors	c19perBeh01, p-value: 0.00052
	MLQ, p-value: 0.00303
R-squared: 0.4809	RSE: 1.183
<b>c19ProSo03</b>	
Best Predictors	age, p-value: 0.00563
R-squared: 0.5432	RSE: 1.156
<b>c19ProSo04</b>	
Best Predictors	consp03, p-value: 0.000146
	c19perBeh02, p-value: 0.000255
R-squared: 0.4152	RSE: 1.075

The table above are the best predictors for pro-social attitudes in similar countries. These 4 models for similar countries have an extremely small p-value which is <2.2e-16, indicating that there is strong evidence against the null hypothesis where the independent variables in the model are strongly associated with the pro-social attitude attributes. The average R-squared value for similar countries is higher than our focus country which indicates that this model has more usefulness than Netherlands, but it may not capture all the underlying relationships in the data. The average RSE for similar countries is lower than Netherlands, indicating these models fit the data better. This suggests that the predictions generated by these models are generally closer to the actual values, enhancing their reliability.

The predictors for similar countries are slightly different than the predictors for Netherlands. Differences in predictors could stem from other various factors such as differences in education systems, government policies and many more. Therefore, the unique characteristics and dynamics of each country could contribute to differences in predictors observed between the Netherlands and other similar countries. Interestingly, we can also see that for Netherlands and its similar countries have common best predictor such as corona personal behavior (c19perBeh01,02,03), indicating that these countries have common understanding of the importance of individual and collective actions in mitigating the impact of the coronavirus pandemic. Similar countries often share similar cultural, socioeconomic, and contextual backgrounds. These similarities can lead to common challenges, priorities, and behaviors among their populations.

When comparing between similar countries and other countries, it seems that the predictors for similar countries gave a more focused and specific understanding of pro-social attitudes. This may be attributed to their similar cultural and socioeconomic backgrounds. While the predictors for other countries capture a wider array of factors, suggesting a more diverse set of influences on pro-social attitudes that may vary across different cultural and societal contexts. It's worth noting that while many predictors for similar countries overlap with those of other countries, the reverse is not true, indicating distinct patterns in the predictors for similar countries.

Call:  
`lm(formula = c19ProSo01 ~ ., data = similar_countries_data)`

Residuals:

	Min	1Q	Median	3Q	Max
	-3.9583	-0.5858	0.0006	0.6597	3.5840

Coefficients: (1 not defined because of singularities)

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.0391693	0.7140324	0.055	0.95627
employstatus_1	-0.0245467	0.1785098	-0.138	0.89067
employstatus_2	-0.0228449	0.1876587	-0.122	0.90314
employstatus_3	-0.1190894	0.1977115	-0.602	0.54713
employstatus_4	0.2856486	0.2142177	1.334	0.18261
employstatus_5	0.0334801	0.2403874	0.139	0.88927
employstatus_6	-0.0787578	0.1954857	-0.403	0.68715
employstatus_7	-0.0621721	0.2192769	-0.284	0.77685
employstatus_8	0.0141530	0.2646614	0.053	0.95737
employstatus_9	-0.0166373	0.2013563	-0.083	0.93417
employstatus_10	-0.1411062	0.3782648	-0.373	0.70923
isoFriends_inPerson	0.0021007	0.0204476	0.103	0.91820
iso0thPpl_inPerson	0.0169974	0.0227685	0.747	0.45558
isoFriends_online	0.0210984	0.0193850	1.088	0.27678
iso0thPpl_online	0.0009694	0.0205672	0.047	0.96242
lone01	0.0539952	0.0612765	0.881	0.37851
lone02	-0.0775124	0.0564502	-1.373	0.17013
lone03	0.0485204	0.0587521	0.826	0.40916
happy	-0.0759035	0.0341300	-2.224	0.02645 *
lifeSat	0.0375941	0.0629342	0.597	0.55045
MLQ	0.0307961	0.0398430	0.773	0.43981
bor01	0.0198225	0.0324544	0.611	0.54153
bor02	0.0194672	0.0335000	0.581	0.56134
bor03	0.1043753	0.0320305	3.259	0.00117 **
consP01	-0.0323420	0.0227951	-1.419	0.15637
consP02	0.0238229	0.0247439	0.963	0.33597
consP03	-0.0130798	0.0204344	-0.640	0.52232
rankOrdLife_1	0.0892775	0.0372767	2.395	0.01687 *
rankOrdLife_2	0.0598799	0.0430581	1.391	0.16474
rankOrdLife_3	0.0081216	0.0442923	0.183	0.85456
rankOrdLife_4	0.0656043	0.0440401	1.490	0.13674
rankOrdLife_5	-0.0084259	0.0425716	-0.198	0.84316
rankOrdLife_6	NA	NA	NA	NA
c19perBeh01	0.1644432	0.0504362	3.260	0.00116 **
c19perBeh02	-0.1545437	0.0611419	-2.528	0.01169 *
c19perBeh03	0.0287734	0.0351153	0.819	0.41282
c19RCA01	0.0699986	0.0298323	2.346	0.01922 *
c19RCA02	0.0104615	0.0444389	0.235	0.81395
c19RCA03	-0.0856204	0.0329662	-2.597	0.00958 **
coronaClose_1	-0.2912803	0.4220190	-0.690	0.49028
coronaClose_2	-0.3773625	0.3272622	-1.153	0.24925
coronaClose_3	0.1035821	0.3415278	0.303	0.76175
coronaClose_4	-0.3141547	0.2618083	-1.200	0.23054
coronaClose_5	-0.2072526	0.2645824	-0.783	0.43369
coronaClose_6	-0.3872678	0.2631603	-1.472	0.14155
gender	0.0331006	0.0842461	0.393	0.69450
age	0.0338501	0.0364018	0.930	0.35273
edu	0.0046545	0.0314164	0.148	0.88226
c19ProSo02	0.2037948	0.0343814	5.927	4.71e-09 ***
c19ProSo03	0.3283178	0.0339533	9.670	< 2e-16 ***
c19ProSo04	0.1114725	0.0385147	2.894	0.00391 **

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 1.134 on 742 degrees of freedom  
Multiple R-squared: 0.4398, Adjusted R-squared: 0.4028  
F-statistic: 11.89 on 49 and 742 DF, p-value: < 2.2e-16

Call:  
`lm(formula = c19ProSo02 ~ ., data = similar_countries_data)`

Residuals:

	Min	1Q	Median	3Q	Max
	-4.9125	-0.5772	0.1140	0.6737	3.1569

Coefficients: (1 not defined because of singularities)

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-1.6970630	0.7423736	-2.286	0.02254 *
employstatus_1	0.2415387	0.1860384	1.298	0.19458
employstatus_2	0.1225991	0.1957429	0.626	0.53129
employstatus_3	0.2759607	0.2060827	1.339	0.18096
employstatus_4	-0.0726773	0.2236608	-0.325	0.74531
employstatus_5	-0.0965199	0.2507851	-0.385	0.70044
employstatus_6	-0.0852892	0.2039572	-0.418	0.67594
employstatus_7	-0.1245652	0.2287481	-0.545	0.58623
employstatus_8	-0.5675887	0.2753462	-2.061	0.03962 *
employstatus_9	0.0489584	0.2100772	0.233	0.81579
employstatus_10	-0.1647970	0.3946512	-0.418	0.67638
isoFriends_inPerson	-0.0005955	0.0213340	-0.028	0.97774
iso0thPpl_inPerson	0.0033055	0.0237640	0.139	0.88941
isoFriends_online	0.0078614	0.0202393	0.388	0.69781
iso0thPpl_online	-0.0072929	0.0214571	-0.340	0.73404
lone01	-0.0523761	0.0639370	-0.819	0.41295
lone02	0.0290792	0.0589621	0.493	0.62203
lone03	0.0690187	0.0612745	1.126	0.26037
happy	0.0722639	0.0356292	2.028	0.04289 *
lifeSat	0.0049902	0.0656775	0.076	0.93945
MLQ	0.1229512	0.0413410	2.974	0.00303 **
bor01	0.0382090	0.0338405	1.129	0.25923
bor02	0.0347968	0.0349367	0.996	0.31958
bor03	-0.0183626	0.0336504	-0.546	0.58544
consP01	-0.0289713	0.0237916	-1.218	0.22372
consP02	-0.0158163	0.0258260	-0.612	0.54045
consP03	0.0462703	0.0212583	2.177	0.02983 *
rankOrdLife_1	0.0736237	0.0389488	1.890	0.05911 .
rankOrdLife_2	0.0651043	0.0449194	1.449	0.14766
rankOrdLife_3	-0.0649832	0.0461516	-1.408	0.15954
rankOrdLife_4	-0.0295222	0.0460049	-0.642	0.52125
rankOrdLife_5	0.0139762	0.0444150	0.315	0.75310
rankOrdLife_6	NA	NA	NA	NA
c19perBeh01	-0.0683903	0.0529384	-1.292	0.19680
c19perBeh02	0.0345934	0.0640535	0.540	0.58931
c19perBeh03	0.1267191	0.0363575	3.485	0.00052 ***
c19RCA01	0.0659674	0.0311467	2.118	0.03451 *
c19RCA02	-0.0333114	0.0463507	-0.719	0.47256
c19RCA03	0.0072744	0.0345501	0.211	0.83330
coronaClose_1	0.7608234	0.4395660	1.731	0.08389 .
coronaClose_2	0.4486797	0.3413559	1.314	0.18912
coronaClose_3	0.3767053	0.3560847	1.058	0.29044
coronaClose_4	-0.0619556	0.2734117	-0.227	0.82080
coronaClose_5	0.1323407	0.2761220	0.479	0.63188
coronaClose_6	0.3426280	0.2746794	1.247	0.21265
gender	-0.1058553	0.0878209	-1.205	0.22845
age	0.0298624	0.0379860	0.786	0.43203
edu	0.0586062	0.0327079	1.792	0.07357 .
c19ProSo01	0.2218444	0.0374265	5.927	4.71e-09 ***
c19ProSo03	0.3587217	0.0352086	10.188	< 2e-16 ***
c19ProSo04	0.0090367	0.0404090	0.224	0.82311

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 1.183 on 742 degrees of freedom  
Multiple R-squared: 0.4809, Adjusted R-squared: 0.4466  
F-statistic: 14.03 on 49 and 742 DF, p-value: < 2.2e-16

```

Call:
lm(formula = c19ProSo03 ~ ., data = similar_countries_data)

Residuals:
    Min      1Q  Median      3Q     Max 
-3.6772 -0.6185  0.0773  0.6380  4.1351 

Coefficients: (1 not defined because of singularities)
                    Estimate Std. Error t value Pr(>|t|)    
(Intercept)       -0.777217   0.726990 -1.069  0.28538  
employstatus_1     -0.143966   0.181814 -0.792  0.42871  
employstatus_2     -0.006939   0.191213 -0.036  0.97106  
employstatus_3      0.107582   0.201464  0.534  0.59350  
employstatus_4     -0.085973   0.218420 -0.394  0.69398  
employstatus_5     -0.131124   0.244894 -0.535  0.59251  
employstatus_6      0.127826   0.199152  0.642  0.52117  
employstatus_7     -0.022637   0.223438 -0.101  0.91933  
employstatus_8     -0.096719   0.269648 -0.359  0.71993  
employstatus_9      0.284635   0.204902  1.389  0.16521  
employstatus_10    0.275993   0.385328  0.716  0.47406  
isoFriends_inPerson 0.019698   0.020822  0.946  0.34446  
iso0thPpl_inPerson -0.031565   0.023179 -1.362  0.17368  
isoFriends_online   0.002890   0.019767  0.146  0.88379  
iso0thPpl_online    0.035544   0.020916  1.699  0.08967 .  
lone01              0.011831   0.062468  0.189  0.84984  
lone02              0.004749   0.057591  0.082  0.93430  
lone03              -0.007435   0.059891 -0.124  0.90124  
happy               -0.006864   0.034891 -0.197  0.84410  
lifeSat              0.069524   0.064090  1.085  0.27837  
MLQ                 -0.039480   0.040588 -0.973  0.33101  
bor01              -0.008633   0.033075 -0.261  0.79416  
bor02              -0.066088   0.034056 -1.941  0.05269 .  
bor03              -0.039982   0.032837 -1.218  0.22376  
consp01             0.008675   0.023256  0.373  0.70924  
consp02             -0.040252   0.025185 -1.598  0.11041  
consp03             0.007301   0.020825  0.351  0.72599  
rankOrdLife_1      -0.044895   0.038093 -1.179  0.23895  
rankOrdLife_2      0.042014   0.043903  0.957  0.33889  
rankOrdLife_3      0.086356   0.045020  1.918  0.05547 .  
rankOrdLife_4      0.042343   0.044914  0.943  0.34611  
rankOrdLife_5      0.056342   0.043329  1.300  0.19389  
rankOrdLife_6          NA      NA      NA      NA      
c19perBeh01         0.009767   0.051756  0.189  0.85038  
c19perBeh02         0.023384   0.062561  0.374  0.70867  
c19perBeh03        -0.060101   0.035728 -1.682  0.09296 .  
c19RCA01            -0.006631   0.030509 -0.217  0.82801  
c19RCA02            -0.048070   0.045247 -1.062  0.28841  
c19RCA03            0.029091   0.033726  0.863  0.38866  
coronaClose_1       0.079212   0.430135  0.184  0.85394  
coronaClose_2       0.063441   0.333747  0.190  0.84929  
coronaClose_3       0.126069   0.347983  0.362  0.71724  
coronaClose_4       0.302702   0.266791  1.135  0.25691  
coronaClose_5       0.199256   0.269603  0.739  0.46010  
coronaClose_6       -0.066365   0.268522 -0.247  0.80486  
gender              0.060651   0.085821  0.707  0.47996  
age                -0.102523   0.036921 -2.777  0.00563 **  
edu                 0.036662   0.031983  1.146  0.25204  
c19ProSo01          0.340865   0.035251  9.670 < 2e-16 ***  
c19ProSo02          0.342129   0.033580 10.188 < 2e-16 ***  
c19ProSo04          0.295339   0.037946  7.783 2.38e-14 ***  
---
```

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 1.156 on 742 degrees of freedom  
 Multiple R-squared: 0.5432, Adjusted R-squared: 0.5131  
 F-statistic: 18.01 on 49 and 742 DF, p-value: < 2.2e-16

```

Call:
lm(formula = c19ProSo04 ~ ., data = similar_countries_data)

Residuals:
    Min      1Q  Median      3Q     Max 
-3.9277 -0.5831  0.0320  0.6716  3.0912 
```

```

Coefficients: (1 not defined because of singularities)
                    Estimate Std. Error t value Pr(>|t|)    
(Intercept)       -0.859842   0.676051 -1.272 0.203821  
employstatus_1     -0.006507   0.169200 -0.038 0.969334  
employstatus_2      0.112915   0.177823  0.635 0.525633  
employstatus_3     -0.012849   0.187443 -0.069 0.945366  
employstatus_4     -0.074712   0.203183 -0.368 0.713196  
employstatus_5      0.152508   0.227782  0.670 0.503362  
employstatus_6     -0.016836   0.185308 -0.091 0.927633  
employstatus_7     -0.029754   0.207847 -0.143 0.886206  
employstatus_8      0.190365   0.250759  0.759 0.448001  
employstatus_9     -0.260697   0.190614 -1.368 0.171829  
employstatus_10    0.508183   0.358081  1.419 0.156266  
isoFriends_inPerson 0.029141   0.019352 -1.506 0.132529  
iso0thPpl_inPerson 0.011695   0.021585  0.542 0.588111  
isoFriends_online   0.034104   0.018346 -1.859 0.063427 .  
iso0thPpl_online    0.025165   0.019472  1.292 0.196636  
lone01              -0.042158   0.058090 -0.726 0.468233  
lone02              0.129690   0.053361  2.430 0.015318 *  
lone03              -0.077013   0.055641 -1.384 0.166745  
happy               0.020503   0.032449  0.632 0.527664  
lifeSat              0.010539   0.059664  0.177 0.859844  
MLQ                 -0.045104   0.037744 -1.195 0.232464  
bor01              -0.012019   0.030766 -0.391 0.696171  
bor02              0.036852   0.031731  1.161 0.245855  
bor03              0.041185   0.030539  1.349 0.177877  
consp01             0.074135   0.021463  3.454 0.000584 ***  
consp02             -0.027743   0.023446 -1.183 0.237082  
consp03             -0.073246   0.019186 -3.818 0.000146 ***  
rankOrdLife_1       0.075357   0.035360  2.131 0.033409 *  
rankOrdLife_2       0.020691   0.040858  0.506 0.612710  
rankOrdLife_3       0.044973   0.041950  1.072 0.284049  
rankOrdLife_4       0.016072   0.041801  0.384 0.700724  
rankOrdLife_5       -0.001187   0.040352 -0.029 0.976544  
rankOrdLife_6          NA      NA      NA      NA      
c19perBeh01         0.075254   0.048067  1.566 0.117867  
c19perBeh02        0.211947   0.057679  3.675 0.000255 ***  
c19perBeh03        0.047212   0.033253  1.420 0.156097  
c19RCA01            -0.010066   0.028378 -0.355 0.722913  
c19RCA02            0.150642   0.041758  3.608 0.000330 ***  
c19RCA03            -0.018791   0.031381 -0.599 0.549491  
coronaClose_1       0.418761   0.399838  1.047 0.295290  
coronaClose_2       0.515352   0.309892  1.663 0.096734 .  
coronaClose_3       0.241158   0.323611  0.745 0.456381  
coronaClose_4       0.336697   0.248084  1.357 0.175135  
coronaClose_5       0.194677   0.250783  0.776 0.437833  
coronaClose_6       0.539852   0.249009  2.168 0.030476 *  
gender              -0.033980   0.079850 -0.426 0.670561  
age                 0.032637   0.034502  0.946 0.344493  
edu                 -0.030831   0.029756 -1.036 0.300494  
c19ProSo01          0.100146   0.034601  2.894 0.003912 **  
c19ProSo02          0.007458   0.033349  0.224 0.823105  
c19ProSo03          0.255564   0.032836  7.783 2.38e-14 ***  
---
```

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 1.075 on 742 degrees of freedom  
 Multiple R-squared: 0.4152, Adjusted R-squared: 0.3766  
 F-statistic: 10.75 on 49 and 742 DF, p-value: < 2.2e-16

## Appendix

I acknowledge the use of ChatGPT (<https://chat.openai.com/>) to generate R codes and refine the language for my own work.

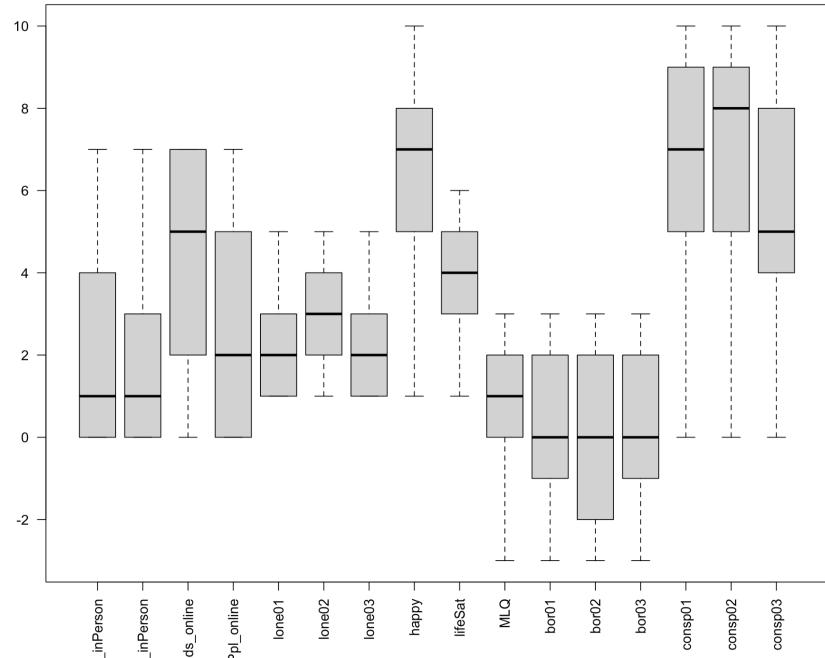


Figure 1.1.1

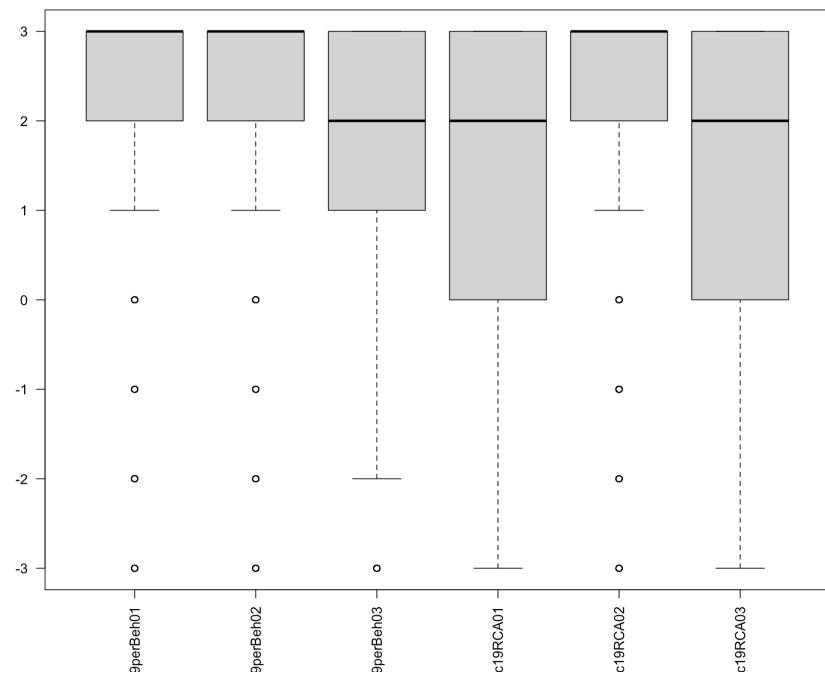


Figure 1.1.2

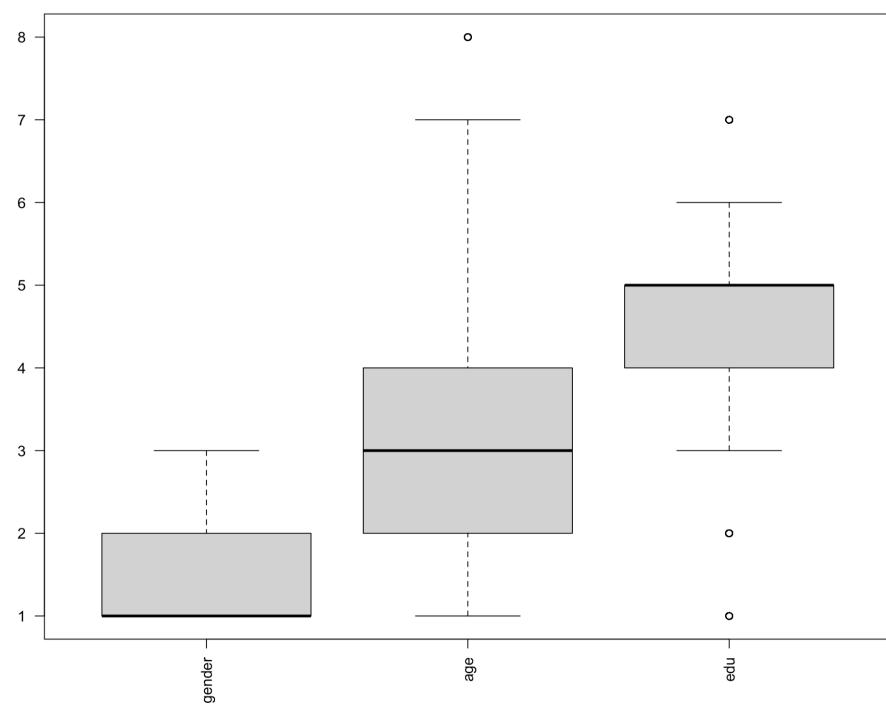


Figure 1.1.3

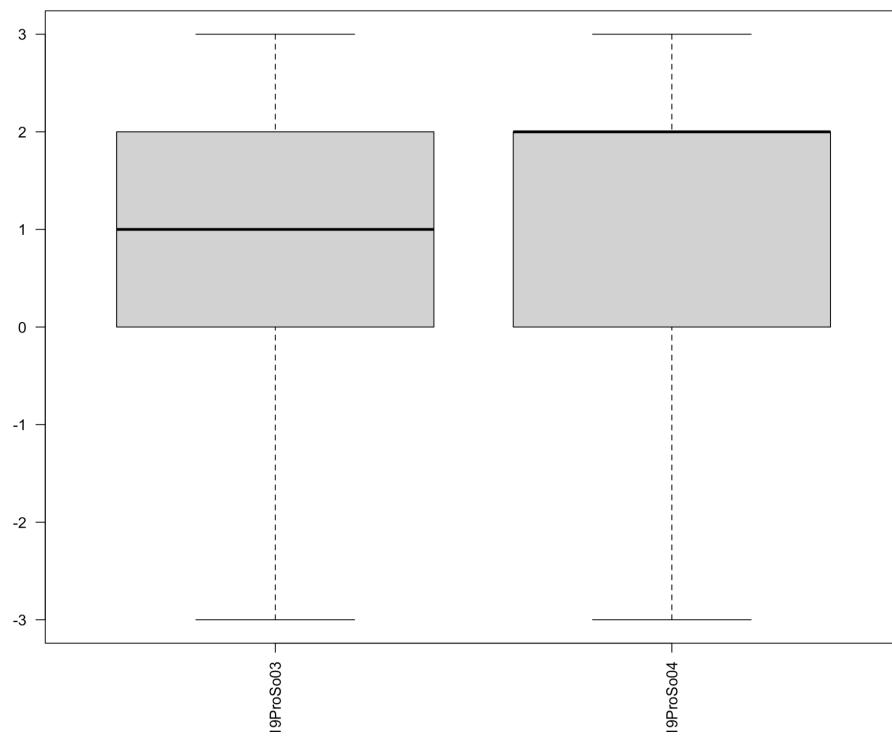


Figure 1.1.4

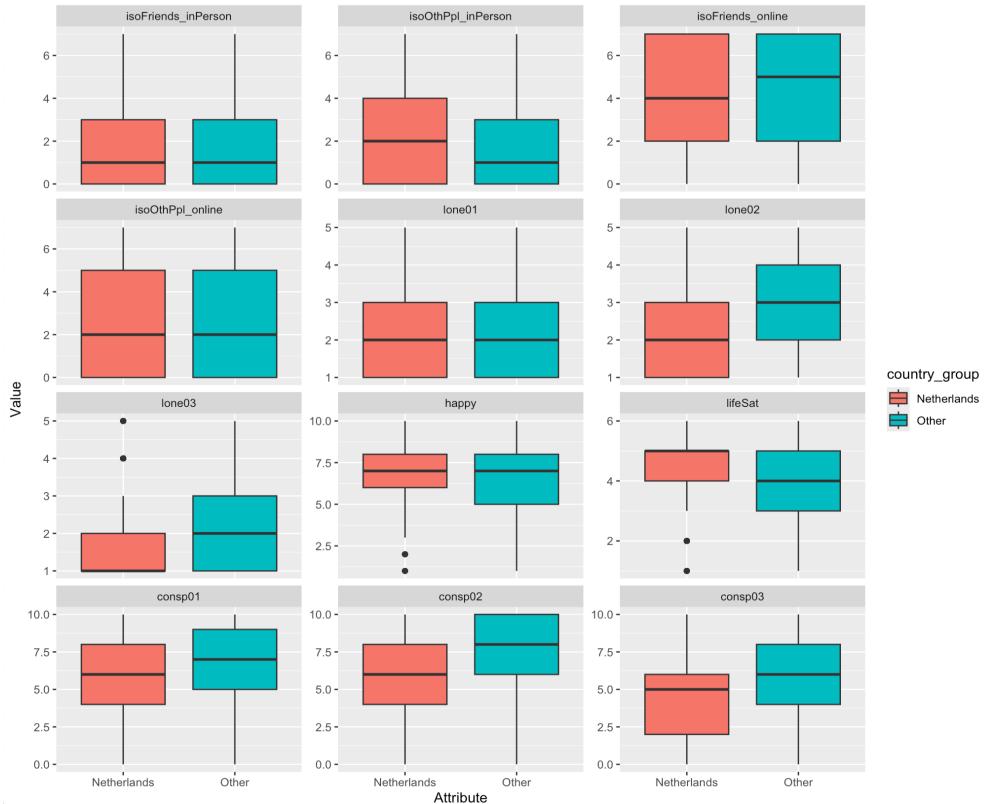


Figure 2.1.1

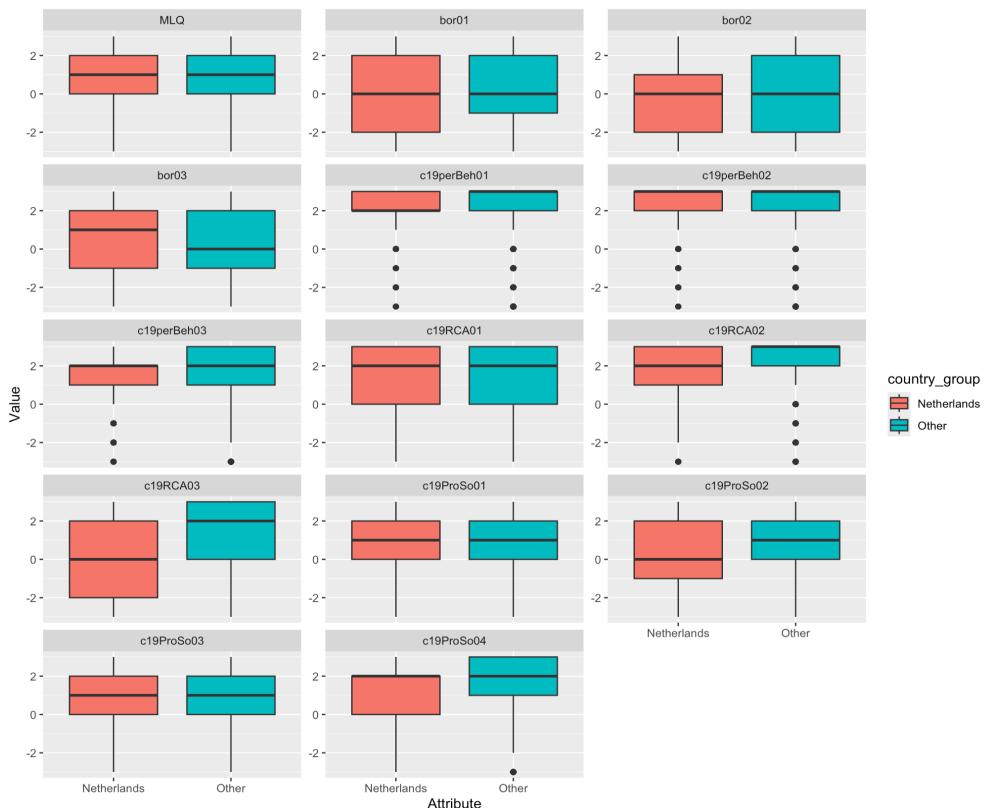


Figure 2.1.2

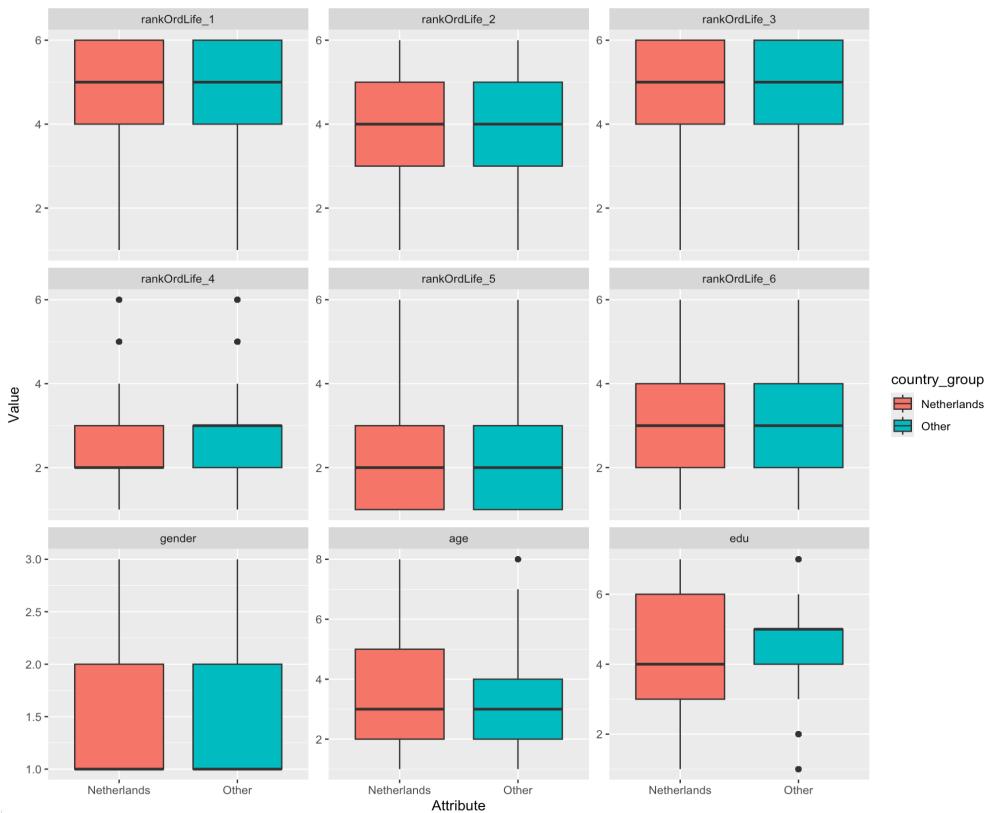


Figure 2.1.3

```

library(dplyr)
library(tidyr)
library(ggplot2)
library(reshape2)
library(cluster)
library(ggrepel)

# reading csv file and create my individual data
rm(list = ls())
set.seed(32909764) # XXXXXXXX = your student ID
cvbase = read.csv("PsyCoronaBaselineExtract.csv")
covid <- cvbase[sample(nrow(cvbase), 40000), ] # 40000 rows

# # Question 1(a)

# getting the rows and columns of the dataset
dim(covid)

# Viewing the structure of the dataset
str(covid)

# Summary statistics for numerical variables
summary(covid)

```

```

# getting all the variables name
colnames(covid)

# checking the data types for each variable
datatypes <- sapply(covid, class)
table(datatypes)

# boxplot for numerical variable
boxplot(covid[11:26], las=2)
boxplot(covid[33:38], las=2)
boxplot(covid[45:47], las=2)
boxplot(covid[51:52], las=2)

# Get the number of unique countries
length(unique(covid$coded_country))

# Create a data frame with frequencies of each country
country_freq <- as.data.frame(table(covid$coded_country))
country_freq <- country_freq[order(-country_freq$Freq), ]

country_freq

# Create a bar plot with sorted countries and modified axis
ggplot(data = country_freq, aes(x = reorder(Var1, Freq), y = Freq)) +
  geom_bar(stat = "identity", fill = "skyblue") +
  labs(x = "Country", y = "Frequency", title = "Distribution of Coded Countries") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

# _____ #
# Question 1(b)

# Columns to replace NA values with 0
columns_to_change_na_to_0 <- c("employstatus_1", "employstatus_2", "employstatus_3",
  "employstatus_4", "employstatus_5", "employstatus_6",
  "employstatus_7", "employstatus_8", "employstatus_9",
  "employstatus_10", "coronaClose_1", "coronaClose_2",
  "coronaClose_3", "coronaClose_4", "coronaClose_5",
  "coronaClose_6")

# Replace NA values with 0 for selected columns
covid <- mutate_at(covid, vars(all_of(columns_to_change_na_to_0)), ~ifelse(is.na(.), 0, .))

# Replace empty strings with NA in the coded_country column
covid <- mutate(covid, coded_country = na_if(coded_country, ""))

# Drop all rows with NA values
covid <- na.omit(covid)

# Define a function to encode categorical variables
encode_categorical <- function(column) {

```

```

# Define encoding mapping
encoding <- c("A" = 1, "B" = 2, "C" = 3, "D" = 4, "E" = 5, "F" = 6)
# Encode the column using the mapping
encoded <- encoding[column]
return(encoded)
}

# Apply the encoding function to each column
covid$rankOrdLife_1 <- encode_categorical(covid$rankOrdLife_1)
covid$rankOrdLife_2 <- encode_categorical(covid$rankOrdLife_2)
covid$rankOrdLife_3 <- encode_categorical(covid$rankOrdLife_3)
covid$rankOrdLife_4 <- encode_categorical(covid$rankOrdLife_4)
covid$rankOrdLife_5 <- encode_categorical(covid$rankOrdLife_5)
covid$rankOrdLife_6 <- encode_categorical(covid$rankOrdLife_6)

# _____ #
# Question 2(a)

# Filter the dataset for records corresponding to the focus country (Netherlands)
netherlands_data <- covid %>% filter(coded_country == "Netherlands")

# Summary statistics for numerical variables in Focus country
summary(netherlands_data)

# Filter the dataset for records corresponding to other countries
other_countries_data <- covid %>% filter(coded_country != "Netherlands")

# Summary statistics for numerical variables in other countries
summary(other_countries_data)

# created a new dataframe and new column to separate between focus country and other countries
two_groups_df <- covid %>%
  mutate(country_group = ifelse(coded_country == "Netherlands", "Netherlands", "Other"))

df_long <- reshape2::melt(two_groups_df, id.vars = "country_group", measure.vars =
  c("isoFriends_inPerson",
    "isoOthPpl_inPerson",
    "isoFriends_online",
    "isoOthPpl_online",
    "lone01", "lone02",
    "lone03", "happy",
    "lifeSat", "consp01",
    "consp02", "consp03"))

# Create side-by-side box plots
ggplot(df_long, aes(x = country_group, y = value, fill = country_group)) +
  geom_boxplot(position = position_dodge(width = 10)) +
  labs(x = "Attribute", y = "Value") +
  facet_wrap(~ variable, scales = "free_y", ncol = 3)

```

```

df_long <- reshape2::melt(two_groups_df, id.vars = "country_group", measure.vars = c("MLQ",
"bor01","bor02",
"bor03", "c19perBeh01",
"c19perBeh02", "c19perBeh03",
"c19RCA01", "c19RCA02",
"c19RCA03", "c19ProSo01",
"c19ProSo02", "c19ProSo03",
"c19ProSo04"))

# Create side-by-side box plots
ggplot(df_long, aes(x = country_group, y = value, fill = country_group)) +
  geom_boxplot(position = position_dodge(width = 10)) +
  labs(x = "Attribute", y = "Value") +
  facet_wrap(~ variable, scales = "free_y", ncol = 3)

df_long <- reshape2::melt(two_groups_df, id.vars = "country_group", measure.vars =
c("rankOrdLife_1", "rankOrdLife_2",
"rankOrdLife_3", "rankOrdLife_4",
"rankOrdLife_5", "rankOrdLife_6",
"gender", "age", "edu"))

# Create side-by-side box plots
ggplot(df_long, aes(x = country_group, y = value, fill = country_group)) +
  geom_boxplot(position = position_dodge(width = 10)) +
  labs(x = "Attribute", y = "Value") +
  facet_wrap(~ variable, scales = "free_y", ncol = 3)

# t-test for isoOthPpl_inPerson
t_test_isoOthPpl_inPerson <- t.test(netherlands_data$isoOthPpl_inPerson,
other_countries_data$isoOthPpl_inPerson)
t_test_isoOthPpl_inPerson

# t-test for consp02
t_test_consp02 <- t.test(netherlands_data$consp02, other_countries_data$consp02)
t_test_consp02

# _____#
# Question 2(b)

# drop coded_country
netherlands_data <- netherlands_data %>% select(-coded_country)

# Perform linear regression for Netherlands
c19ProSo01_lm <- lm(c19ProSo01 ~ ., data = netherlands_data)
c19ProSo02_lm <- lm(c19ProSo02 ~ ., data = netherlands_data)
c19ProSo03_lm <- lm(c19ProSo03 ~ ., data = netherlands_data)
c19ProSo04_lm <- lm(c19ProSo04 ~ ., data = netherlands_data)

# Print summary of the regression model for Netherlands
summary(c19ProSo01_lm)
summary(c19ProSo02_lm)

```

```

summary(c19ProSo03_lm)
summary(c19ProSo04_lm)

# Question 2(c)

# drop coded_country
other_countries_data <- other_countries_data %>% select(-coded_country)

# Perform linear regression for other countries
c19ProSo01_lm_oth <- lm(c19ProSo01 ~ ., data = other_countries_data)
c19ProSo02_lm_oth <- lm(c19ProSo02 ~ ., data = other_countries_data)
c19ProSo03_lm_oth <- lm(c19ProSo03 ~ ., data = other_countries_data)
c19ProSo04_lm_oth <- lm(c19ProSo04 ~ ., data = other_countries_data)

# Print summary of the regression model for other countries
summary(c19ProSo01_lm_oth)
summary(c19ProSo02_lm_oth)
summary(c19ProSo03_lm_oth)
summary(c19ProSo04_lm_oth)

# _____ #

# Question 3(a)

# Read indicator datasets
GDP_per_capita_indicator <- read.csv("recent_GDP_capita.csv")
Urban_population_indicator <- read.csv("recent_urban_population.csv")
Rural_population_indicator <- read.csv("recent_rural_population_percent.csv")
birth_rate_indicator <- read.csv("recent_fctb_birth_rate.csv")
death_rate_indicator <- read.csv("recent_fctb_death_rate.csv")

# Merge the indicator datasets by the column "country"
indicator_df <- merge(GDP_per_capita_indicator, Urban_population_indicator, by = "country")
indicator_df <- merge(indicator_df, Rural_population_indicator, by = "country")
indicator_df <- merge(indicator_df, birth_rate_indicator, by = "country")
indicator_df <- merge(indicator_df, death_rate_indicator, by = "country")

# Drop rows with missing values
indicator_df <- na.omit(indicator_df)

# Keep the "country" column as an identifier
countries <- indicator_df$country

# Remove the "country" column before scaling the data
indicator_df <- indicator_df[, -which(names(indicator_df) == "country")]

# Standardize the data
scaled_data <- scale(indicator_df)

# Determine the number of clusters (k)
# elbow method:
wss <- numeric(10)

```

```

for (i in 1:10) {
  wss[i] <- sum(kmeans(scaled_data, centers = i)$withinss)
}
plot(1:10, wss, type = "b", xlab = "Number of clusters", ylab = "Within-cluster sum of squares", main = "Elbow Method")

# silhouette method:
sil_width <- sapply(2:10, function(k) {
  kmeans_result <- kmeans(scaled_data, centers = k)
  silhouette_width <- silhouette(kmeans_result$cluster, dist(scaled_data))
  mean(silhouette_width[, "sil_width"])
})

# Plot silhouette widths for different numbers of clusters
plot(2:10, sil_width, type = "b", xlab = "Number of clusters", ylab = "Mean Silhouette Width", main = "Silhouette Method")

# Find the optimal number of clusters with the maximum silhouette width
optimal_clusters <- which.max(sil_width) + 1
cat("Optimal number of clusters based on silhouette method:", optimal_clusters, "\n")

# Perform k-means clustering with the optimal number of clusters
optimal_kmeans <- kmeans(scaled_data, centers = optimal_clusters)

# Add cluster assignments to the indicator dataframe
indicator_df$cluster <- as.factor(optimal_kmeans$cluster)

# Reattach the "countries" vector to the dataframe
indicator_df <- cbind(country = countries, indicator_df)

# Convert cluster column to numeric
indicator_df$cluster <- as.numeric(as.character(indicator_df$cluster))

# Find the cluster assignment of the Netherlands
netherlands_cluster <- as.integer(optimal_kmeans$cluster[indicator_df$country == "Netherlands"])

# Filter dataframe to include only countries in the same cluster as the Netherlands
same_cluster_df <- indicator_df[indicator_df$cluster == netherlands_cluster, ]

# Calculate distances between each country in the same cluster and the Netherlands using Euclidean distance
country_distances <- apply(same_cluster_df[, -1], 1, function(x) {
  sqrt(sum((x - same_cluster_df[same_cluster_df$country == "Netherlands", -1])^2))
})

# Add distances as a new column to the dataframe
same_cluster_df$distance_to_netherlands <- country_distances

# Sort the dataframe based on distances
sorted_same_cluster_df <- same_cluster_df[order(same_cluster_df$distance_to_netherlands), ]

```

```

# Extract rows 2 to 7 from the sorted dataframe (excluding the Netherlands)
closest_countries <- sorted_same_cluster_df$country[2:7]

# Print the closest countries
closest_countries

# Add cluster assignments to the scaled data
scaled_data_with_clusters <- cbind(scaled_data, cluster = optimal_kmeans$cluster)

# Convert the data to a data frame
scaled_data_with_clusters_df <- as.data.frame(scaled_data_with_clusters)

# Plot the cluster plot
ggplot(scaled_data_with_clusters_df, aes(x = latest.value_NY.GDP.PCAP.CD, y =
latest.value_SP.URB.TOTL.IN.ZS, color = factor(cluster))) +
  geom_point() +
  scale_color_discrete(name = "Cluster") +
  labs(x = "GDP per Capita", y = "Urban Population (% of total)", title = "Cluster Plot") +
  theme_minimal()

# _____ #

# Question 3(b)

# Filter the covid dataset with the selected similar countries
similar_countries_data <- covid[covid$coded_country %in% closest_countries, ]

# drop coded_country
similar_countries_data <- similar_countries_data %>% select(-coded_country)

# Linear regression for similar countries
lm_similar_c19ProSo01 <- lm(c19ProSo01 ~ ., data = similar_countries_data)
lm_similar_c19ProSo02 <- lm(c19ProSo02 ~ ., data = similar_countries_data)
lm_similar_c19ProSo03 <- lm(c19ProSo03 ~ ., data = similar_countries_data)
lm_similar_c19ProSo04 <- lm(c19ProSo04 ~ ., data = similar_countries_data)

# Summary of linear regression for similar countries
summary(lm_similar_c19ProSo01)
summary(lm_similar_c19ProSo02)
summary(lm_similar_c19ProSo03)
summary(lm_similar_c19ProSo04)

# Filter the covid dataset without the selected similar countries
other_countries_data <- covid[!covid$coded_country %in% closest_countries, ]

# drop coded_country
other_countries_data <- other_countries_data %>% select(-coded_country)

# Linear regression for other countries
lm_other_c19ProSo01 <- lm(c19ProSo01 ~ ., data = other_countries_data)
lm_other_c19ProSo02 <- lm(c19ProSo02 ~ ., data = other_countries_data)
lm_other_c19ProSo03 <- lm(c19ProSo03 ~ ., data = other_countries_data)

```

```
lm_other_c19ProSo04 <- lm(c19ProSo04 ~ ., data = other_countries_data)

# Summary of linear regression for similar countries
summary(lm_other_c19ProSo01)
summary(lm_other_c19ProSo02)
summary(lm_other_c19ProSo03)
summary(lm_other_c19ProSo04)
```