

Reproducibility and Literate Programming

Rachel

2022-10-19

For this portfolio I will be using data obtained from the World Bank's DataBank on Gender Statistics [^1] to demonstrate how one might use literate programming to produce code that can easily be recreated by anyone reading this report. We first load the packages that we'll be making use of in our analysis [^1]: <https://databank.worldbank.org/reports.aspx?source=gender-statistics#>

Since there are approximately 1000 possible metrics and 265 countries to consider for a period of >20 years, I will only be focusing on the employment data for 2020, which I have downloaded into a CSV file directly from the website. To load this into our environment we can run the following (as long as the data is in the same directory as the markdown file)

```
data <- read_csv("Raw_Gender_Data.csv", show_col_types = FALSE)
```

If we look at a summary of the first two columns we can see they are just denoting the year, which is constant and thus not useful, thus we remove it. The fourth column contains a unique series code and the third column contains a longer description of the metric. Since this is the same information, we can remove the less concise third column. If we want to retrieve it, the data comes with a CSV file containing the series name for each series column.

```
data <- data[, -c(1,2,3)]
```

Then we might want to take a brief visual inspection of our data (for simplicity we just look at the first 15 rows and 5 countries):

```
knitr::kable(data[1:15, 1:6], booktabs = TRUE) %>%  
  kable_styling(font_size = 8, latex_options = "scale_down")
```

Series Code	United Kingdom [GBR]	Afghanistan [AFG]	Albania [ALB]	Algeria [DZA]	American Samoa [ASM]
SG.GET.JOBS.EQ	1	1	1	1	..
SG.NGT.WORK.EQ	1	0	1	0	..
SG.DNG.WORK.DN.EQ	1	0	1	0	..
SG.IND.WORK.EQ	1	0	1	1	..
SH.HIV.ARTC.FE.ZS	..	9	53	87	..
SH.HIV.ARTC.MA.ZS	..	9	45	80	..
SE.PRM.TENR
SE.PRM.TENR.FE
SE.PRM.TENR.MA
SP.ADO.TFRT	11.183	57.509	19.4332	9.3594	..
SH.HIV.PMTC.ZS	..	10	..	34	..
SH.STA.BRTC.ZS
SH.DTH.COMM.ZS
SH.DTH.COMM.0004.ZS
SH.DTH.COMM.0004.FE.ZS

Number of Missing Values	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
--------------------------	---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----

We can see that there are many missing values, although here they are represented by an ellipses. We can then change these to NA, as this is what R recognises as missing values, and then see how many missing values there are in each row

```
data <- data %>% na_if("..")
na_freq <- tibble("Number of Missing Values" = 1:max(rowSums(is.na(data))), "Frequency" = tabulate(rowSums(is.na(data)), 100))
na_freq <- na_freq[!(na_freq$Frequency == 0),]
knitr::kable(t(na_freq))%>%
  kable_styling(font_size = 8, latex_options = "scale_down")
```