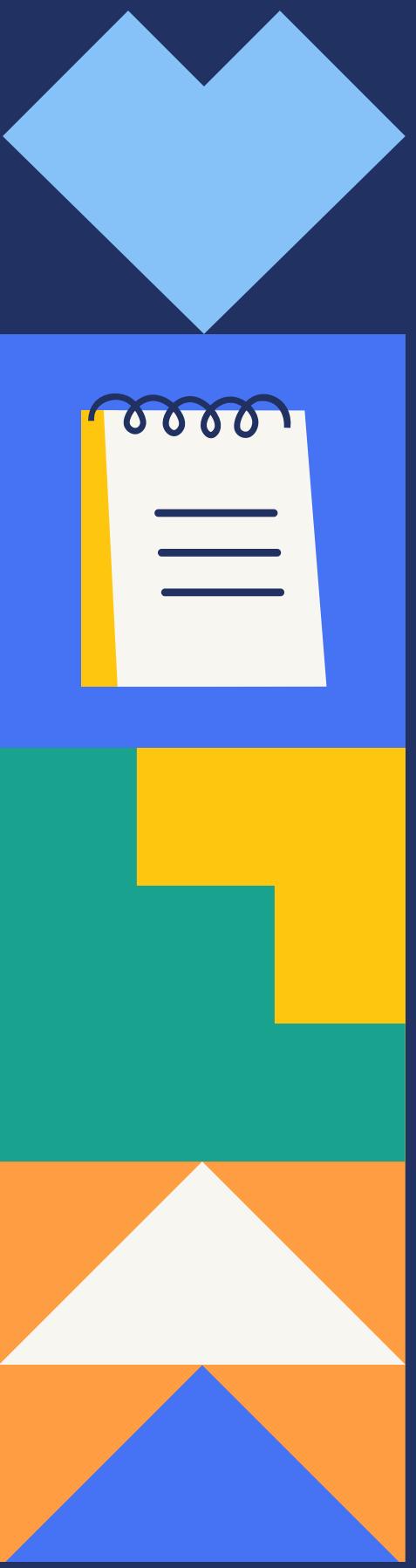


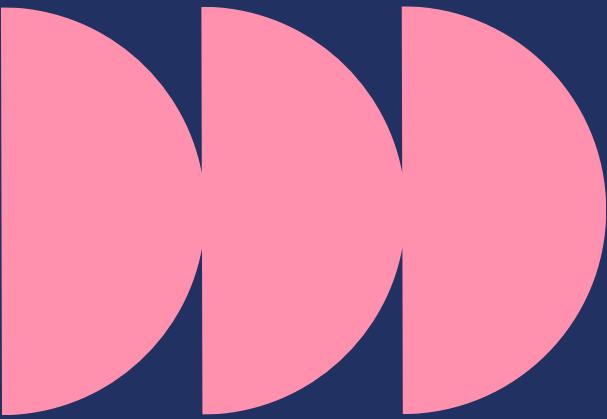
Diagnosing and Curing Data Problems (Using R Statistical Software!)

Alexia Samiotis
PhD Candidate (Clinical Neuropsychology)
Monash Epworth Rehabilitation Research Centre





Outline



Diagnosing Data Problems

- Visualising
- Descriptives
- Missing Data
- Outliers

R Crash Course



Curing Data Problems

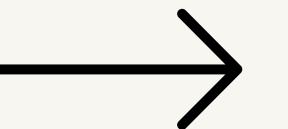
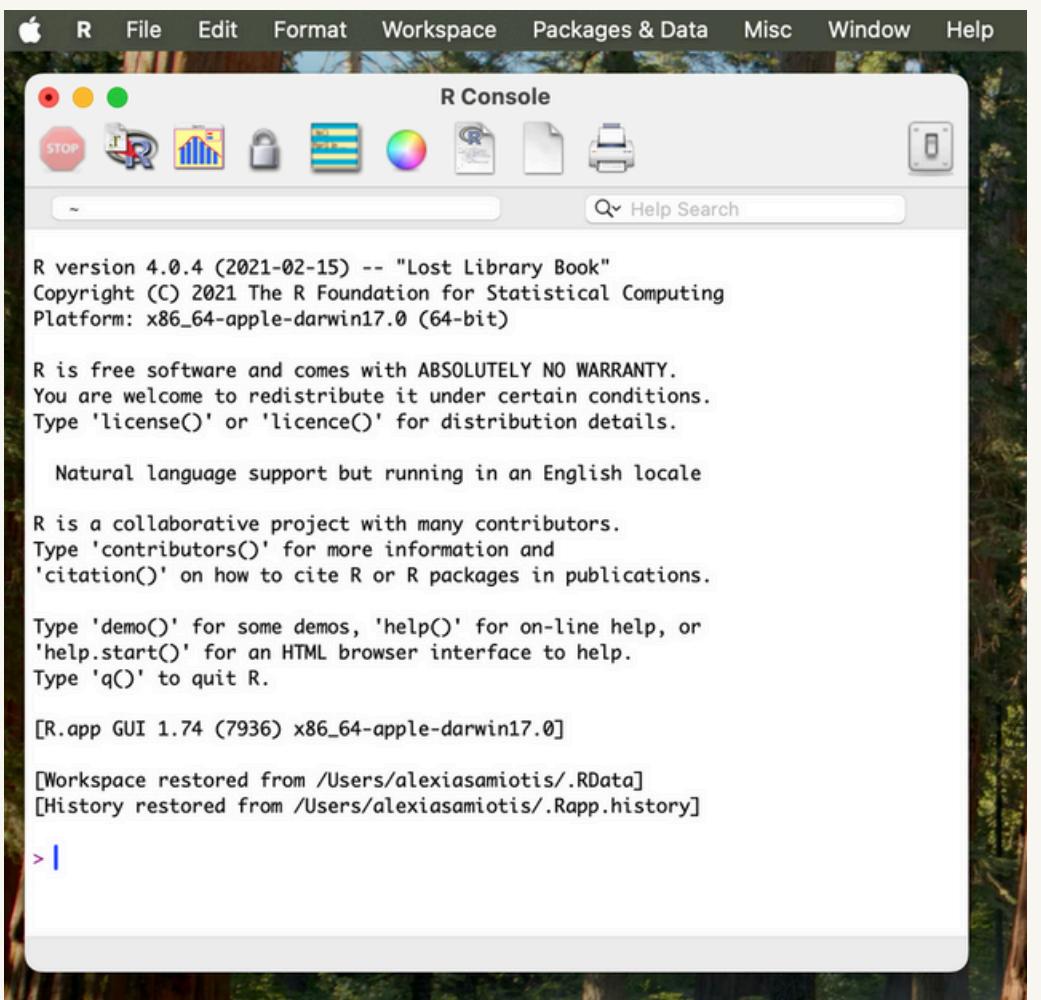
- Handling missing data
- Handling Outliers
- Visualising relationships

Link to code



R Crash Course

R - computer programming language

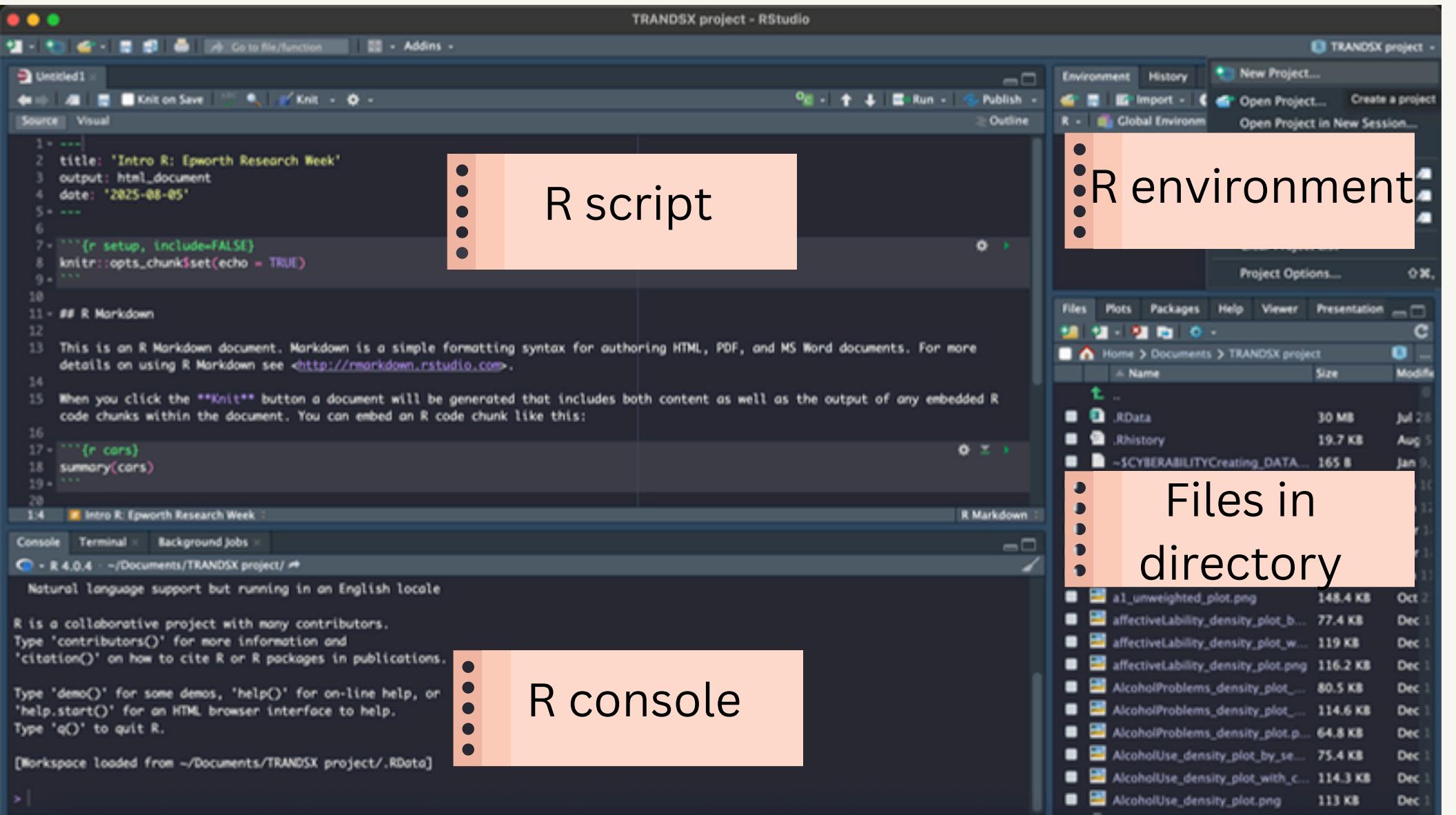


A screenshot of a DataCamp article page. The title is "How to Install R on Windows, Mac OS X, and Ubuntu Tutorial". The page content includes a green header with the DataCamp logo, a brief description, and a timestamp "DataCamp / Mar 11, 2020". The main content area features a person sitting at a laptop with code snippets overlaid.

Downloading R and R Studio to Mac if needed

R Crash Course

R Studio - An application that provides an interface that looks the same on Mac/Windows and is more user friendly



Downloading to R Studio to Mac if needed



How to Install R on Windows, Mac OS X, and Ubuntu Tutorial

This is a beginner guide that is designed to save yourself a headache and valuable time if you deci...

DataCamp / Mar 11, 2020

View on DataCamp

Read on Medium

Share on LinkedIn

Share on Facebook

Share on Twitter

Share on Email

Copy link

Report abuse

Save to my profile

Save to my library

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

Save to my bookmarks

Save to my reading list

Save to my notes

Save to my tasks

Save to my calendar

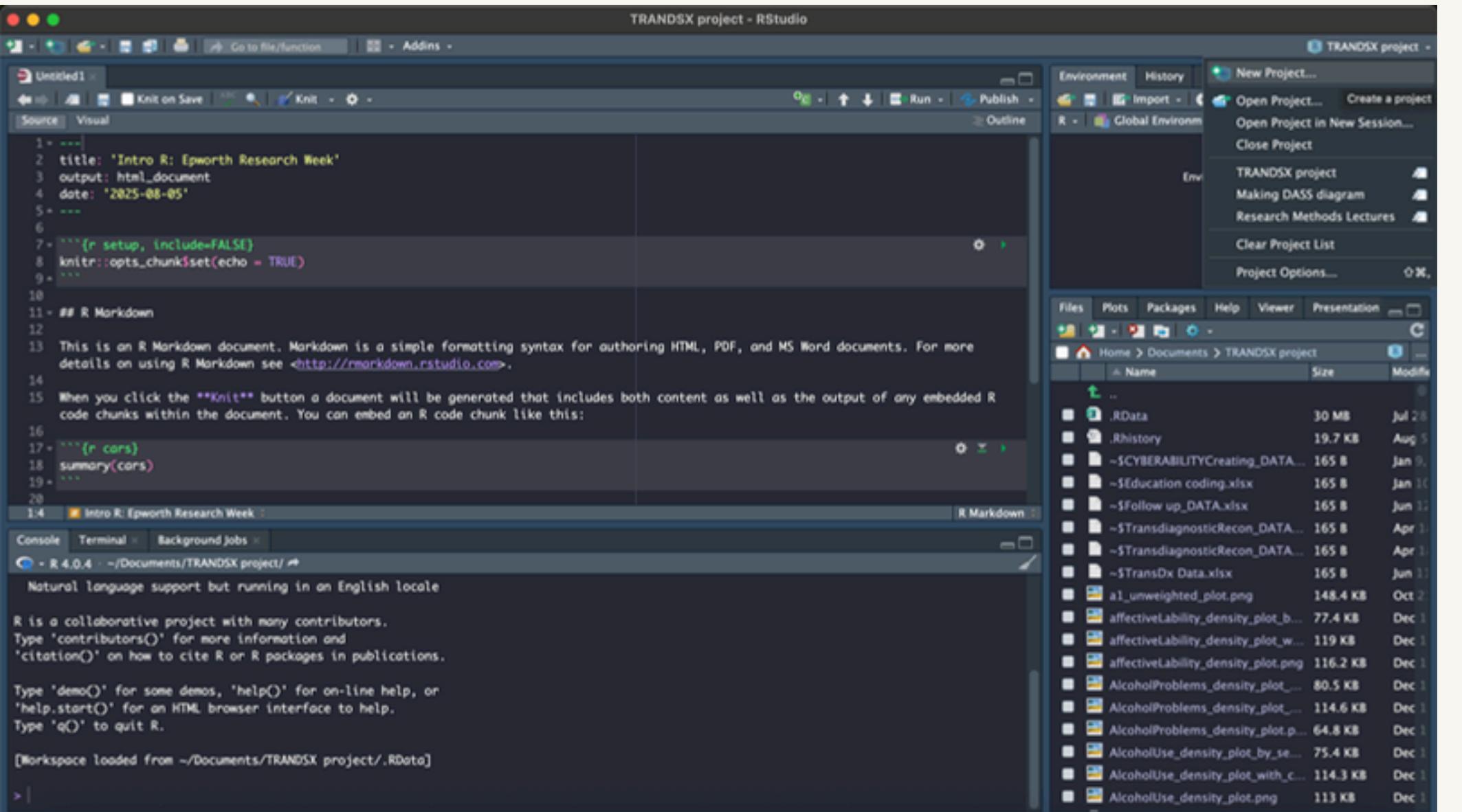
Save to my bookmarks

Save to my reading list

R Crash Course

Setting Up

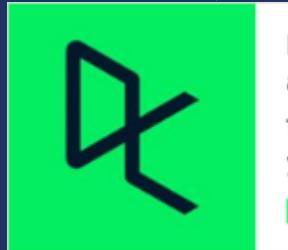
1. New Project



The screenshot shows the RStudio interface with the following details:

- Title Bar:** TRANSDX project - RStudio
- Source Editor:** Contains R Markdown code for a document titled "Intro R: Epworth Research Week". The code includes setup chunks for knitr and R code blocks.
- Console:** Displays the R environment, version R 4.0.4, and a message about natural language support.
- File Explorer:** Shows the project structure with files like .RData, .Rhistory, and various XLSX and PNG files related to "SCYBERABILITYCreating_DATA..." and "TransdiagnosticRecon_DATA...".
- Project Menu:** Opened, showing options like "New Project...", "Open Project...", and "Project Options...".

Downloading to R Studio to Mac if needed



How to Install R on Windows, Mac OS X, and Ubuntu Tutorial

This is a beginner guide that is designed to save yourself a headache and valuable time if you deci...

DataCamp / Mar 11, 2020



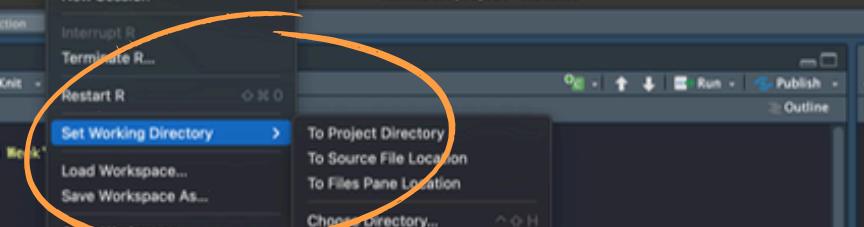
R Crash Course



Setting Up

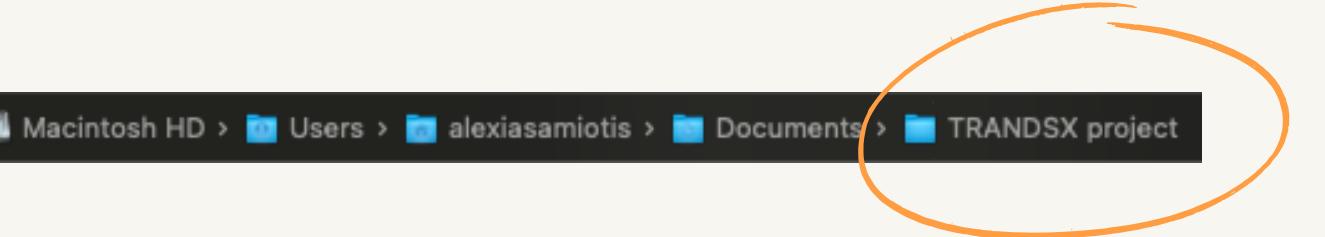
1. New Project

2. Set Working Directory

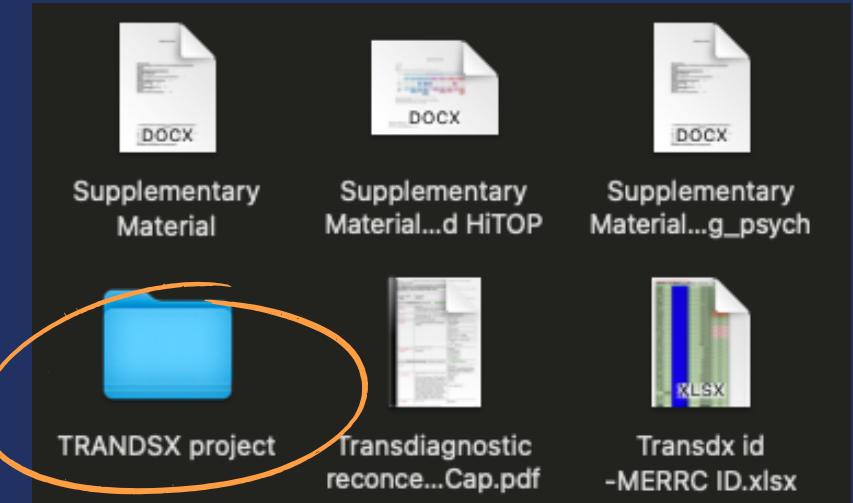


```
Macintosh HD > Users > alexiasamiotis > Documents > TRANDSX project
```

The screenshot shows the RStudio interface. In the top-left corner, the 'Session' menu is open, and the 'Set Working Directory' option is highlighted with an orange circle. The main workspace shows an R Markdown file named 'Untitled1.Rmd'. The code in the file includes R code chunks and a note about R Markdown. The bottom-left corner shows the R console output, which includes the path to the working directory: 'R - R 4.0.4 - ~/Documents/TRANDSX project/'. The bottom-right corner shows the file browser with a list of files in the 'TRANDSX project' folder.



Create a new folder on your computer itself (NOT OneDrive or iCloud)



Set this as your project directory. Mine is in my Documents folder.

R Crash Course

Setting Up

1. New Project

2. Set Working Directory

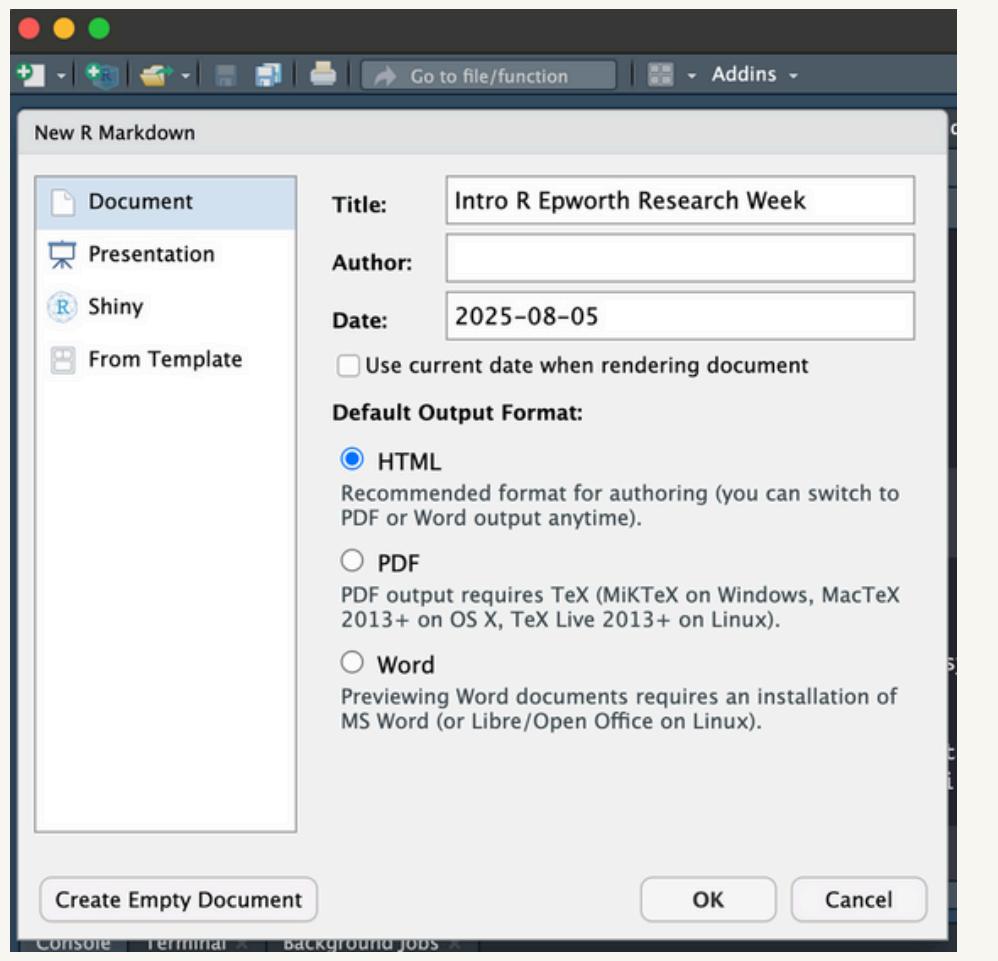
3. Open new Markdown file

The screenshot shows the RStudio interface. On the left, the 'File' menu is open, and the 'R Markdown' option is highlighted with a red circle. The main workspace shows an R Markdown document titled 'Epworth Research Week'. The code chunk contains the following R code:

```
dev<-FALSE>
set(echo = TRUE)

## This is a simple R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more information see http://rmarkdown.rstudio.com.
```

The 'Knit' button is visible in the top right of the R Markdown editor.



Now we create a new file within R. I prefer Markdown.

R Markdown

Turn your analyses into high quality documents, reports, presentations and dashboards with R Markdown. Use a productive notebook interface to weave together narrative text and code to produce elegantly formatted...

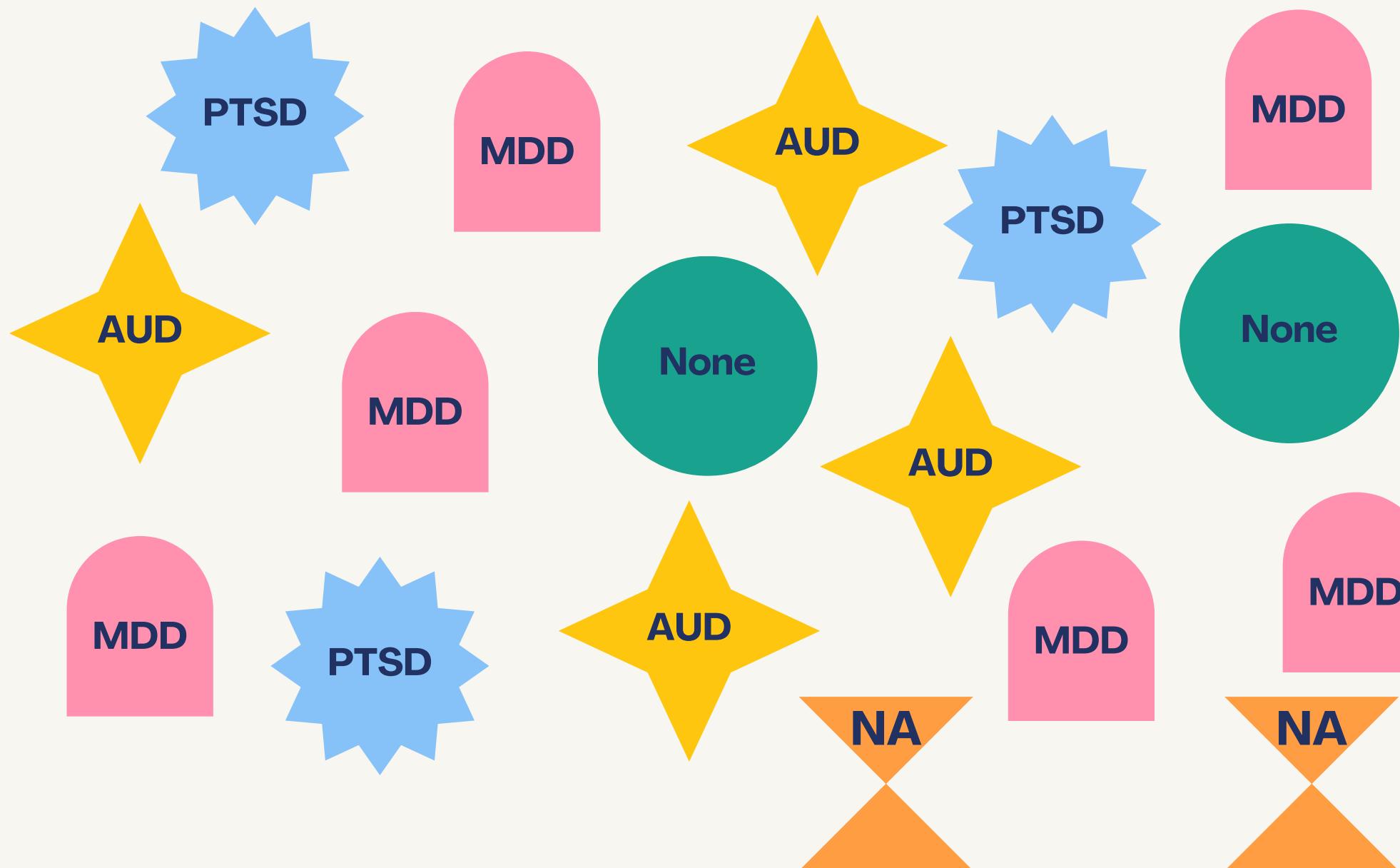
rstudio.com

You can knit it at the end to a nice looking file, like the one I made for today.



Diagnosing Problems: Discreet Data

Psychologists interviewed individuals with TBI
and these are the diagnoses that were provided:



How can we first organise and
present the data in R?

Associations between transdiagnostic psychopathology dimensions and cognitive functioning

after traumatic brain injury: An application of the HiTOP-TBI model

Alexia Samiotis^{*1,2}, Jai Carmichael^{1,2}, Jao-Yue Carmintati^{1,2}, Amelia J Hicks¹, Jennie Ponsford^{1,2}, Kate Rachel Gould^{1,2}, Gershon Spitz^{1,2,3}

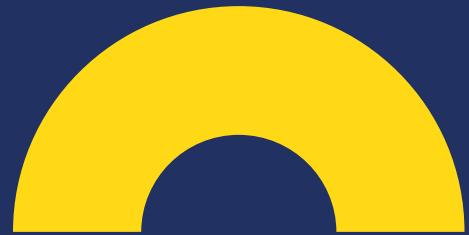
Importing data

Associations between transdiagnostic psychopathology dimensions and cognitive functioning

after traumatic brain injury: An application of the HiTOP-TBI model

Alexia Samiotis^{*1,2}, Jai Carmichael^{1,2}, Jao-Yue Carmintati^{1,2}, Amelia J Hicks¹, Jennie

Ponsford^{1,2}, Kate Rachel Gould^{1,2}, Gershon Spitz^{1,2,3}



Discrete Data

Represent the frequency of different results in the data.

So how do we get this to R Studio?

You should first have a worksheet where you entered the data

result or
data value



DSM
Diagnosis

MDD

frequency



Number of Individuals
with TBI

5

AUD

12

-
-
-
-
-

1. Installing packages

```
17 - ## Installing Packages
18 These include functions that you will need to execute basic commands like importing your data or creating basic plots, exporting data. You only have to install these ONCE, not every time you open R (thankfully!). In fact, if you try and install them more than once you will effectively just be updating the package (which is fine just not necessary).
19
20 - ```{r packages}
21 install.packages("readxl") # needed to import data from excel sheets
22 install.packages("tidyverse") # meta-package with other packages I often use
23 install.packages("tidyverse") # for data tidying
24 install.packages("lubridate") # for data wrangling
25 install.packages("writexl") # for exporting to excel
26 install.packages("ggplot2") # visualisation, plots
27 install.packages("psych") # psych stats stuff :(
28 install.packages("openxlsx") # used sometimes when readxl doesn't work
29 -
30 ````
```

I like my R Studio interface with these colours, you can change yours to have a white background if you like

Tools > Options > Environment > General and choose a colour theme from the "Colour theme" dropdown menu

2. Opening packages

```
31 - ## Load Relevant Libraries
32 You can think of installing packages like downloading an app, once it's done it will be there. This doesn't automatically mean the app will open, you still have to click on it. In R, you have to tell the package to open by loading it in the environment. You do have to load packages EVERY time you use R Studio. If you stay in the same session and don't exit R, you don't have to load them.
33
34 - ```{r libraries, echo=FALSE}
35 library(readxl)
36 library(tidyverse)
37 library(tidyr)
38 library(lubridate)
39 library(writexl)
40 library(ggplot2)
41 library(psych)
42 library(openxlsx)
43 ````
```



Link to code

3.Importing and naming your data

```
45 ## Import the dataset
46 This is unpublished data that I have collected as part of my PhD looking at mental health after Traumatic Brain Injury. As you can
see there >1000 variables, many of which will be included in the analyses for my studies. Today we will only be focusing on a few of
these. I will show you how to call on the variables you would like to use.
47
48 ````{r transdiagnostic data}
49 df <- read_excel("HiTOP_followup_data.xlsx") # I called this df but can call it something more meaningful if you like.
50 view(df) # lets have a look at what our imported dataset looks like
51 ````
```



Naming: what would you like to name the dataframe? name <-

I like my R Studio interface with these colours, you can change yours to have a white background if you like

Tools > Options > Environment >
General and choose a colour theme from the "Colour theme" dropdown menu

Link to code

[https://github.com/AlexiaSam/R-Intro/blob/main_Epworth Research Week Intro.md](https://github.com/AlexiaSam/R-Intro/blob/main_Epworth%20Research%20Week%20Intro.md)



3.Importing and naming your data

This is the function to import the data (found in the `readxl` package we loaded earlier)

Type the name of your excel file (must be verbatim)

```
45 ## Import the dataset  
46 This is unpublished data that I have collected as part of my PhD looking at mental health after Traumatic Brain Injury. As you can see there >1000 variables, many of which will be included in the analyses for my studies. Today we will only be focusing on a few of these. I will show you how to call on the variables you would like to use.  
47  
48 #'r transdiagnostic data}  
49 df <- read_excel("HiTOP_followup_data.xlsx") # I called this df but can call it something more meaningful if you like.  
50 view(df) # lets have a look at what our imported dataset looks like  
51 #'
```

Naming: what would you like to name the dataframe? `name <-`

I like my R Studio interface with these colours, you can change yours to have a white background if you like

Tools > Options > Environment > General and choose a colour theme from the "Colour theme" dropdown menu

[Link to code](#)



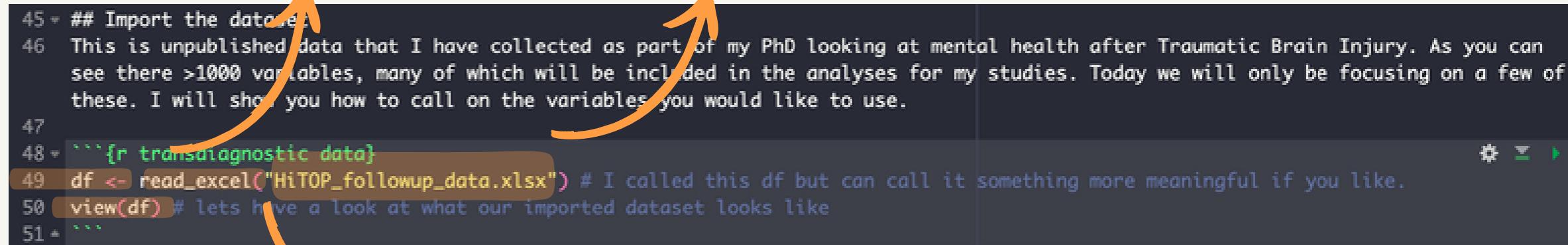
[https://github.com/AlexiaSam/R-Intro/blob/main_Epworth Research Week Intro.md](https://github.com/AlexiaSam/R-Intro/blob/main_Epworth%20Research%20Week%20Intro.md)

3.Importing and naming your data

This is the function to import the data (found in the `readxl` package we loaded earlier)

Type the name of your excel file (must be verbatim)

```
45 ## Import the dataset
46 This is unpublished data that I have collected as part of my PhD looking at mental health after Traumatic Brain Injury. As you can see there >1000 variables, many of which will be included in the analyses for my studies. Today we will only be focusing on a few of these. I will show you how to call on the variables you would like to use.
47
48 ```{r transdiagnostic data}
49 df <- read_excel("HiTOP_followup_data.xlsx") # I called this df but can call it something more meaningful if you like.
50 view(df) # lets have a look at what our imported dataset looks like
51 ```

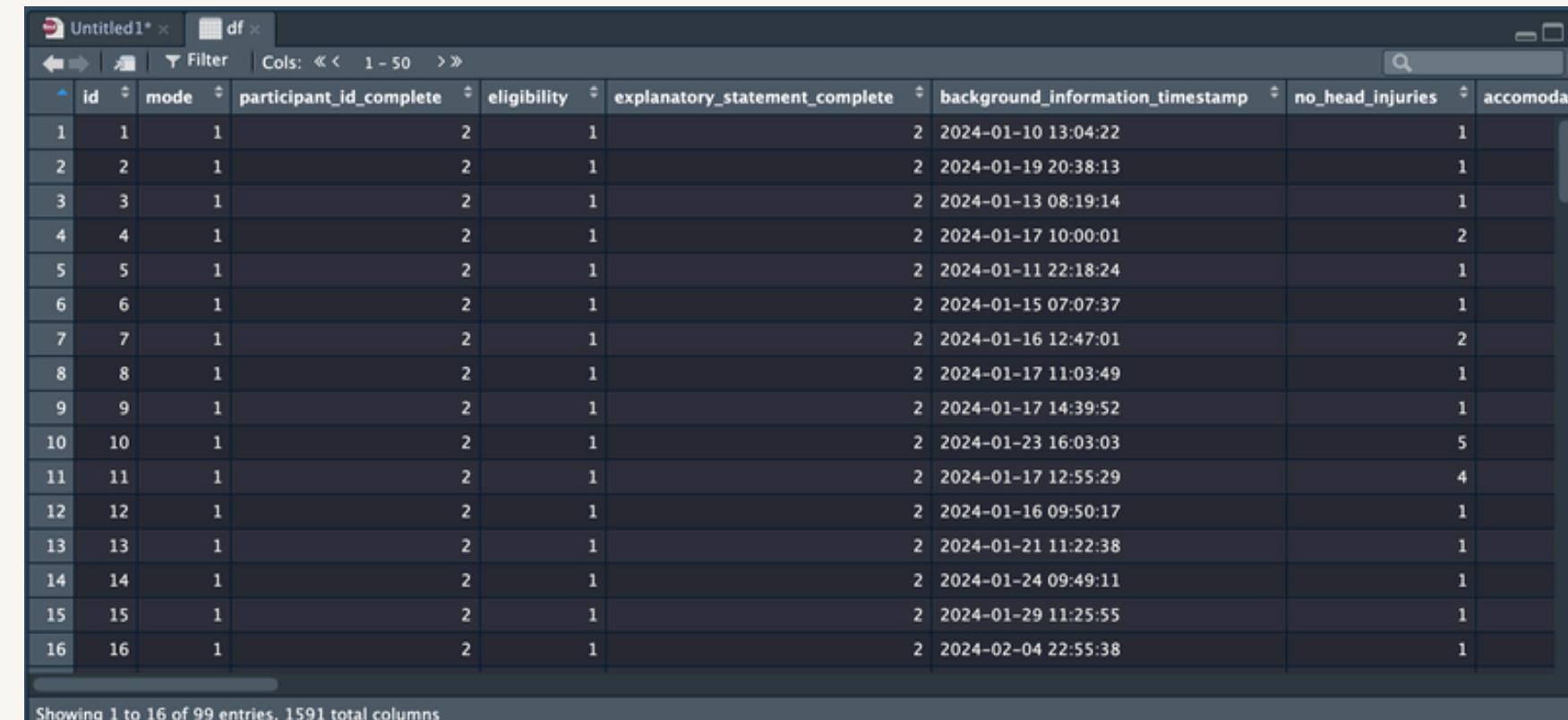

```

Let's have a look at our dataframe

Naming: what would you like to name the dataframe? `name <-`

I like my R Studio interface with these colours, you can change yours to have a white background if you like

Tools > Options > Environment > General and choose a colour theme from the "Colour theme" dropdown menu



	id	mode	participant_id_complete	eligibility	explanatory_statement_complete	background_information_timestamp	no_head_injuries	accomodat	...
1	1	1		2	1	2 2024-01-10 13:04:22		1	
2	2	1		2	1	2 2024-01-19 20:38:13		1	
3	3	1		2	1	2 2024-01-13 08:19:14		1	
4	4	1		2	1	2 2024-01-17 10:00:01		2	
5	5	1		2	1	2 2024-01-11 22:18:24		1	
6	6	1		2	1	2 2024-01-15 07:07:37		1	
7	7	1		2	1	2 2024-01-16 12:47:01		2	
8	8	1		2	1	2 2024-01-17 11:03:49		1	
9	9	1		2	1	2 2024-01-17 14:39:52		1	
10	10	1		2	1	2 2024-01-23 16:03:03		5	
11	11	1		2	1	2 2024-01-17 12:55:29		4	
12	12	1		2	1	2 2024-01-16 09:50:17		1	
13	13	1		2	1	2 2024-01-21 11:22:38		1	
14	14	1		2	1	2 2024-01-24 09:49:11		1	
15	15	1		2	1	2 2024-01-29 11:25:55		1	
16	16	1		2	1	2 2024-02-04 22:55:38		1	

4. Calling on specific variables

\$ is used to call the variable

```
53 - ## Calling on specific variables
54 I'm sure you don't want to go searching through 1000 variables, so instead you can call on the variable you need and then work with
it. You can also perform functions on a vector of variables you call on (1). Alternatively, you can filter a subset of variables that
you would like to work with (2). I'll show you both.
55
56 - ``{r calling on variables}
57 df$mini_mdd_current # what does our Major Depressive Disorder variable contain? Let's have a quick look by printing it here.
58 df$mini_alcohol_current # AUD diagnoses
59 df$mini_ptsd_current # PTSD diagnoses
60 - ````
```

Which dataframe contains the variable you would like to call on? Mine is `df`

Which variable do you want to call?
You must type verbatim the name
that is in your dataframe.

What does the output look like?

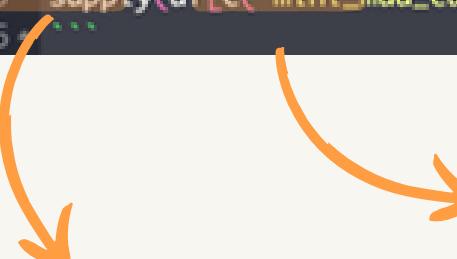
How should we present our diagnostic data?

There are some things to consider (Tufte, 2001):

1. Show the data
2. Induce the reader to think about the data being presented
3. Avoid distorting the data
4. Present many numbers with minimum ink
5. Make large datasets coherent
6. Encourage the reader to compare different bits of data

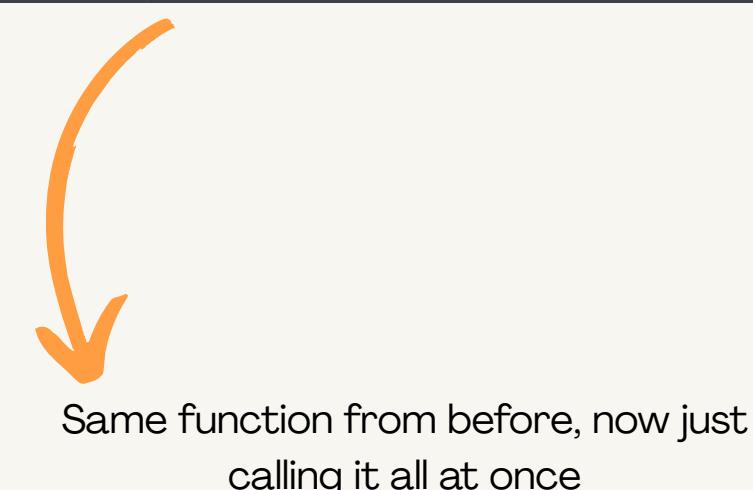
5 people have MDD
13 people have AUD
8 people have PTSD

```
73 + ### 1. Show the data
74
75 Check what they contain by calling on each variable individually (option 1)
76 This gives you the frequency of each diagnosis. For this dataset, 1 = NO diagnosis, 2 = YES diagnosis
77 + ````{r}
78 table(df$mini_mdd_current)
79 table(df$mini_alcohol_current)
80 table(df$mini_ptsd_current)
81 +
82
83 OR Calling on variables altogether as a vector (option 2)
84 + ````{r}
85 sapply(df[c("mini_mdd_current", "mini_alcohol_current", "mini_ptsd_current")], table)
86 + ````{r}
```



Variables called on in `vector` =
`c("variable", "variable")`

Apply the function to all
variables in the vector



```
Console Terminal × Background Jobs ×
R 4.0.4 ~/Documents/TRANDSX project/
> table(df$mini_mdd_current)
1 2
66 5
> table(df$mini_alcohol_current)
1 2
65 13
> table(df$mini_ptsd_current)
1 2
70 8
> sapply(df[c("mini_mdd_current", "mini_alcohol_current", "mini_ptsd_current")], table)
mini_mdd_current mini_alcohol_current mini_ptsd_current
1 66 65 70
2 5 13 8
> |
```

How should we present our diagnostic data?

There are some things to consider (Tufte, 2001):

1. Show the data
2. Induce the reader to think about the data being presented
3. Avoid distorting the data
4. Present many numbers with minimum ink
5. Make large datasets coherent
6. Encourage the reader to compare different bits of data

[Link to code](#)

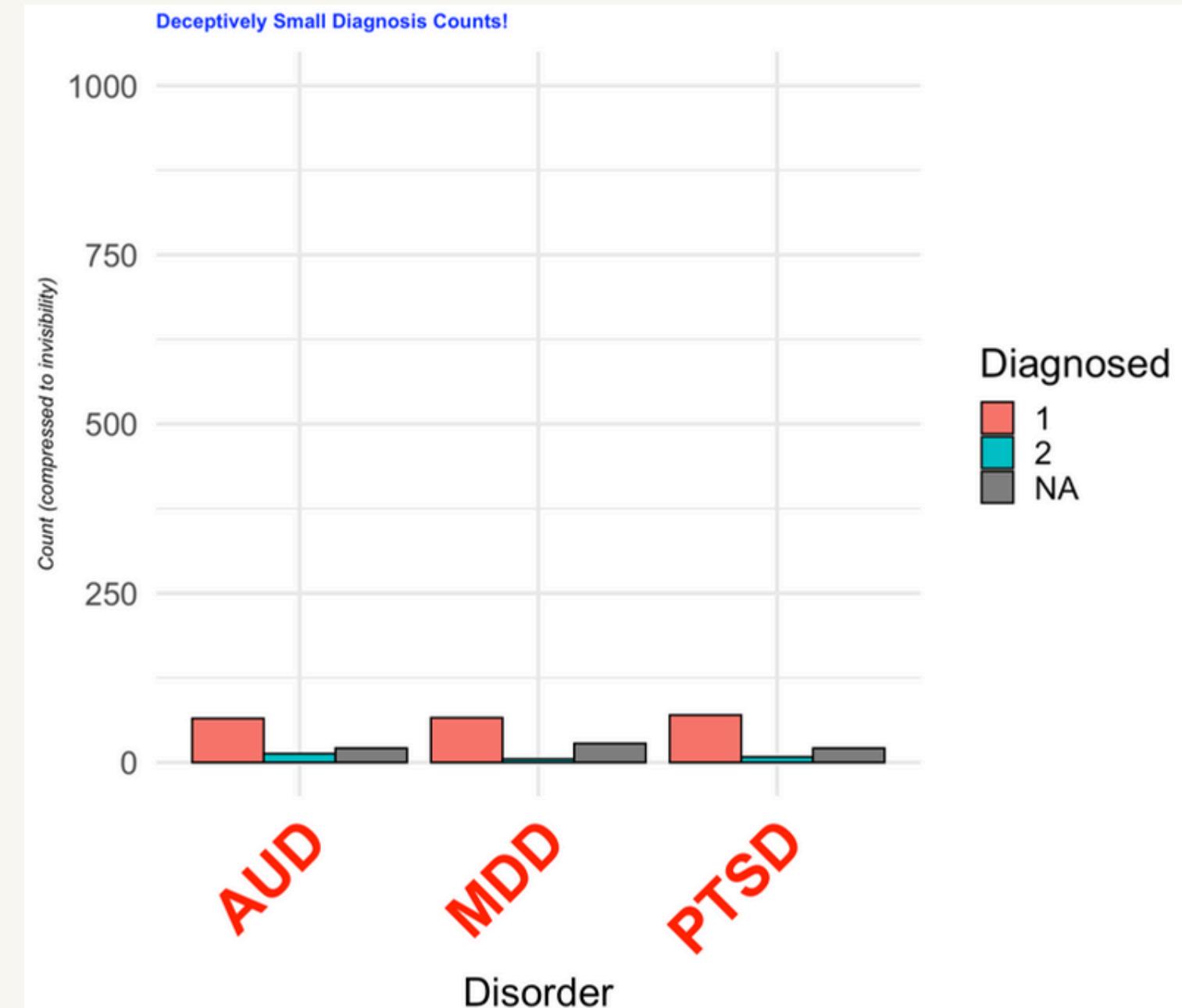


How should we present our diagnostic data?

There are some things to consider (Tufte, 2001):

1. Show the data
2. Induce the reader to think about the data being presented
3. **Avoid distorting the data**
4. Present many numbers with minimum ink
5. Make large datasets coherent
6. Encourage the reader to compare different bits of data

Anything wrong with this?



How should we present our diagnostic data?

There are some things to consider (Tufte, 2001):

1. Show the data
2. Induce the reader to think about the data being presented
3. **Avoid distorting the data**
4. Present many numbers with minimum ink
5. Make large datasets coherent
6. Encourage the reader to compare different bits of data

Y-axis inflation	- By setting the y-limit to 10,00 (likely way above actual values), you minimise visual differences, making everything look tiny.
Perception distortion	- Viewers may falsely assume diagnoses are rare or not worth attention.
Undermines comparison	- Tiny differences become invisible, reducing insight.



How should we present our diagnostic data?

There are some things to consider (Tufte, 2001):

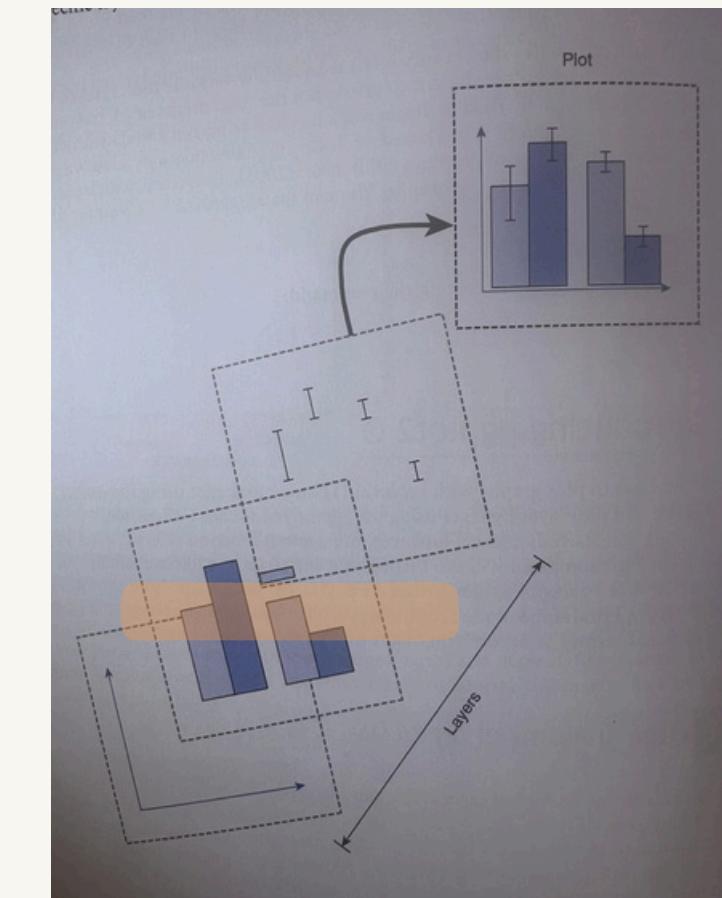
1. Show the data
2. Induce the reader to think about the data being presented
3. Avoid distorting the data
4. Present many numbers with minimum ink
5. Make large datasets coherent
6. Encourage the reader to compare different bits of data

What does this look like in R?

```
93 2. Induce the reader to think about the frequency of diagnoses across each category using ggplot to create a bar chart.  
94  
95 ````{r}  
96 df %>% # calling on our df  
97   select(MDD = mini_mdd_current, # elect the variables to plot  
98         AUD = mini_alcohol_current,  
99         PTSD = mini_ptsd_current) %>%  
100 pivot_longer(everything(), names_to = "Diagnosis", values_to = "Status") %>% # changing our data to long format as is  
better for plotting  
101 group_by(Diagnosis, Status) %>% # specify which groups to plot and how  
102 summarise(Count = n(), .groups = 'drop') %>%  
103 ggplot(aes(x = Diagnosis, y = Count, fill = factor(Status))) + # layer to label things nicely and tell the graph where to  
overlay the data  
104   geom_bar(stat = "identity", position = "dodge") +  
105   labs(fill = "Diagnosed", title = "Current Diagnoses per MINI") + # add a title layer  
106   theme_minimal()  
107  
108 ````
```

Layer to tell ggplot what to fill
the axes with

geom_bar - tells ggplot to add
a bar plot layer
Layer for the axis labels and
titles

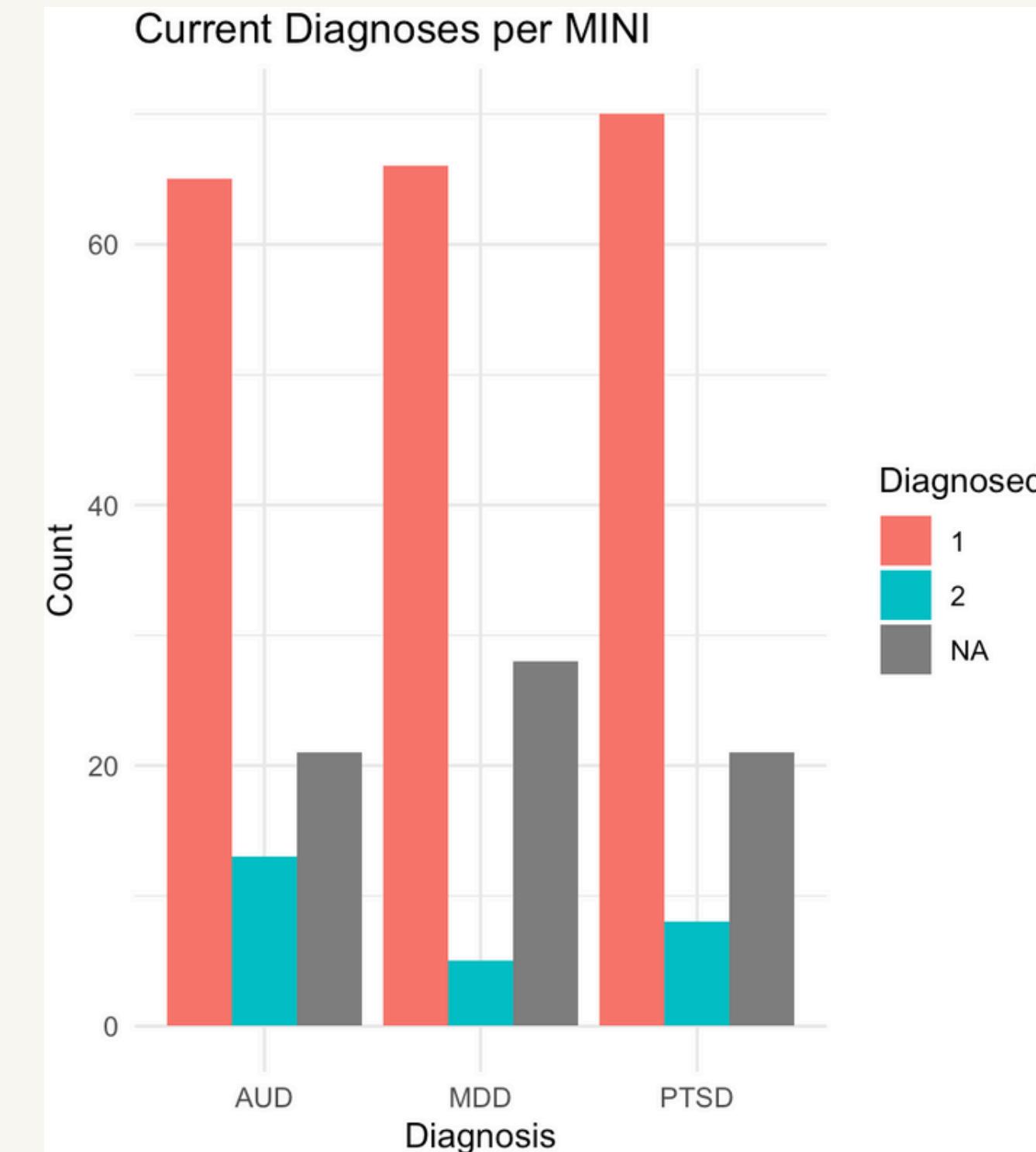


How should we present our diagnostic data?

There are some things to consider (Tufte, 2001):

1. Show the data ✓
2. Induce the reader to think about the data being presented ✓
3. Avoid distorting the data ✓
4. Present many numbers with minimum ink ✓
5. Make large datasets coherent ✓
6. Encourage the reader to compare different bits of data ✓

What does this look like in R?
Pretty!



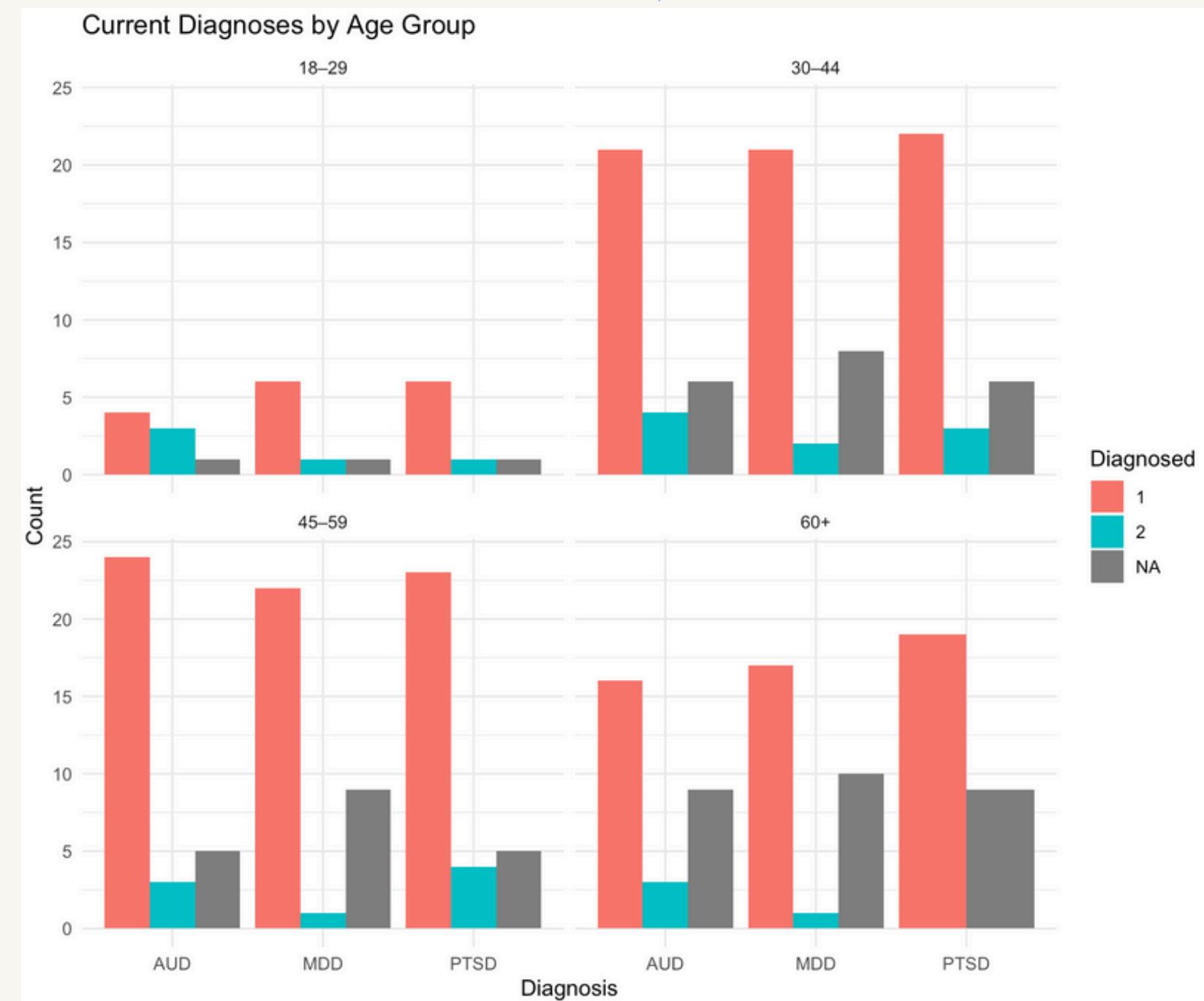
Bonus: stratify by other variables

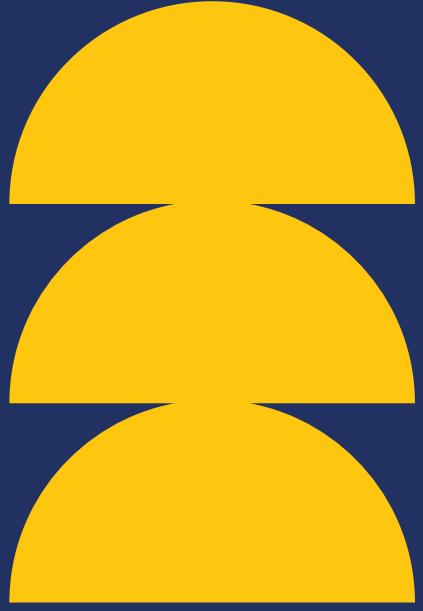
There are some things to consider (Tufte, 2001):

1. Show the data ✓
2. Induce the reader to think about the data being presented ✓
3. Avoid distorting the data ✓
4. Present many numbers with minimum ink ✓
5. Make large datasets coherent ✓
6. Encourage the reader to compare different bits of data ✓

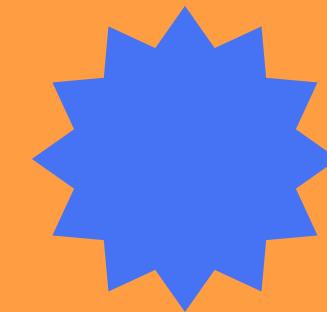
Pretty!

```
148 ## BONUS: stratify the diagnoses by age bracket and make look pretty :)
149
150 ``{r}
151 # first need to specify the age groups
152 df$age_group <- cut(df$age, breaks = c(18, 30, 45, 60, 100),
153   labels = c("18-29", "30-44", "45-59", "60+"),
154   right = FALSE)
155 df %>%
156   select(age_group,
157     MDD = mini_mdd_current,
158     AUD = mini_alcohol_current,
159     PTSD = mini_ptsd_current) %>%
160   pivot_longer(cols = c(MDD, AUD, PTSD), names_to = "Diagnosis", values_to = "Status") %>%
161   group_by(age_group, Diagnosis, Status) %>%
162   summarise(Count = n(), .groups = "drop") %>%
163   ggplot(aes(x = Diagnosis, y = Count, fill = factor(Status))) +
164   geom_bar(stat = "identity", position = "dodge") +
165   facet_wrap(~age_group) +
166   labs(title = "Current Diagnoses by Age Group", fill = "Diagnosed") +
167   theme_minimal()
168 ``
```





Problem 1: Missing Data



Frequency

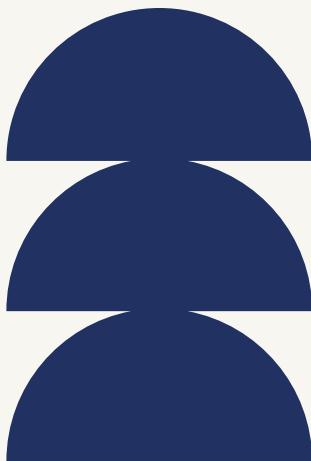
we have plotted already
(NAs)



Solution

different options - for diagnostic data, we will be transparent about missingness and filter missing data for analysis

Problem 1: Missing Data



filter the data
!= NOT in R
is.na = anything with an NA
value

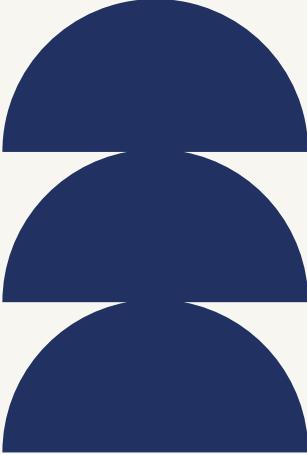
Filter the data to only include individuals without NAs across any of our diagnostic variables

Solution

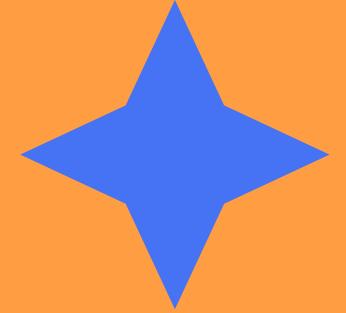
different options - for diagnostic data, we will be transparent about missingness and filter missing data for analysis

```
171 ## Problem 1: Handling Missing Diagnoses in the data
172 #While there are many methods to handle missing data, for today's purpose we will
173 #missingness (achieved above).
174
175 #### Solution: We will then filter the data so it is clean prior to analysis.
176
177 ``{r}
178 df_clean <- df %>%
179   filter(!is.na(mini_mdd_current) & # filter the dataset to only include values WITHOUT NAs
180         !is.na(mini_alcohol_current) &
181         !is.na(mini_ptsd_current))
182
183 # Total rows before cleaning
184 total_rows <- nrow(df)
185
186 # Total rows after removing missing
187 clean_rows <- nrow(df_clean)
188
189 # Rows with any missing in key diagnosis variables
190 missing_rows <- total_rows - clean_rows
```

Define a variable for the missingness. Simple maths formula: Total rows - clean rows (which you calculate in the two steps above!)



Problem 1: Missing Data



Solution

Plot data without NAs
and add footnote about
missingness

```
210 ## Plot clean dataset
211
212 ````{r}
213 ggplot(plot_data, aes(x = Diagnosis, y = Count, fill = factor(Status))) +
214   geom_bar(stat = "identity", position = "dodge") +
215   geom_text(aes(label = Label),
216             position = position_dodge(width = 0.9),
217             vjust = -0.5, size = 5) +
218   scale_fill_manual(values = c("0" = "#999999", "1" = "#0072B2"),
219                     labels = c("Not Diagnosed", "Diagnosed")) +
220   labs(
221     title = "Current Diagnoses (Cleaned Data)",
222     subtitle = paste("Missing data removed: ", missing_rows, " of ", total_rows, " cases (",
223                     round(100 * missing_rows / total_rows, 1), "%)", sep = ""),
224     x = "Disorder",
225     y = "Count",
226     fill = "Diagnosed"
227   ) +
228   theme_minimal(base_size = 16)
229
230 ````
```

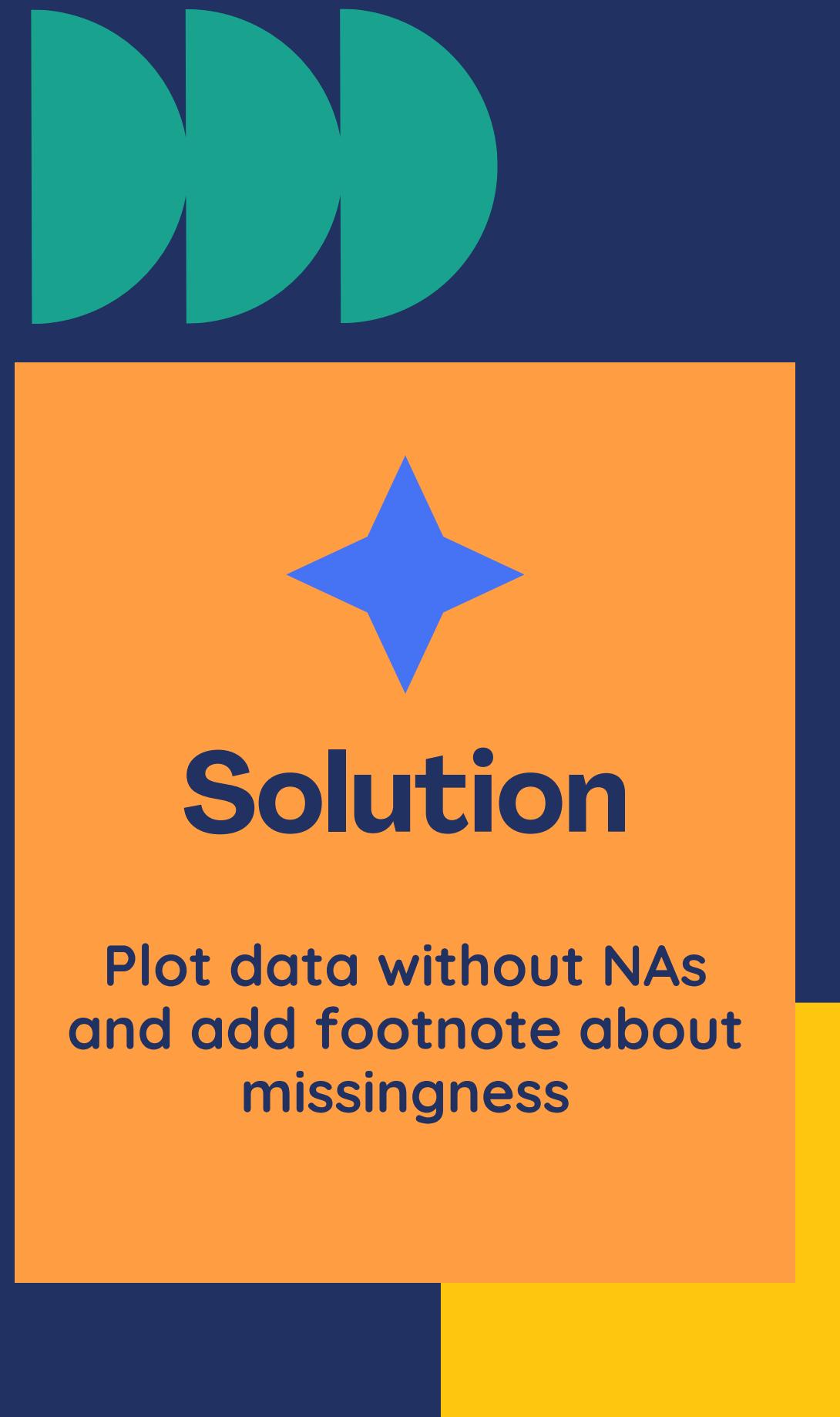
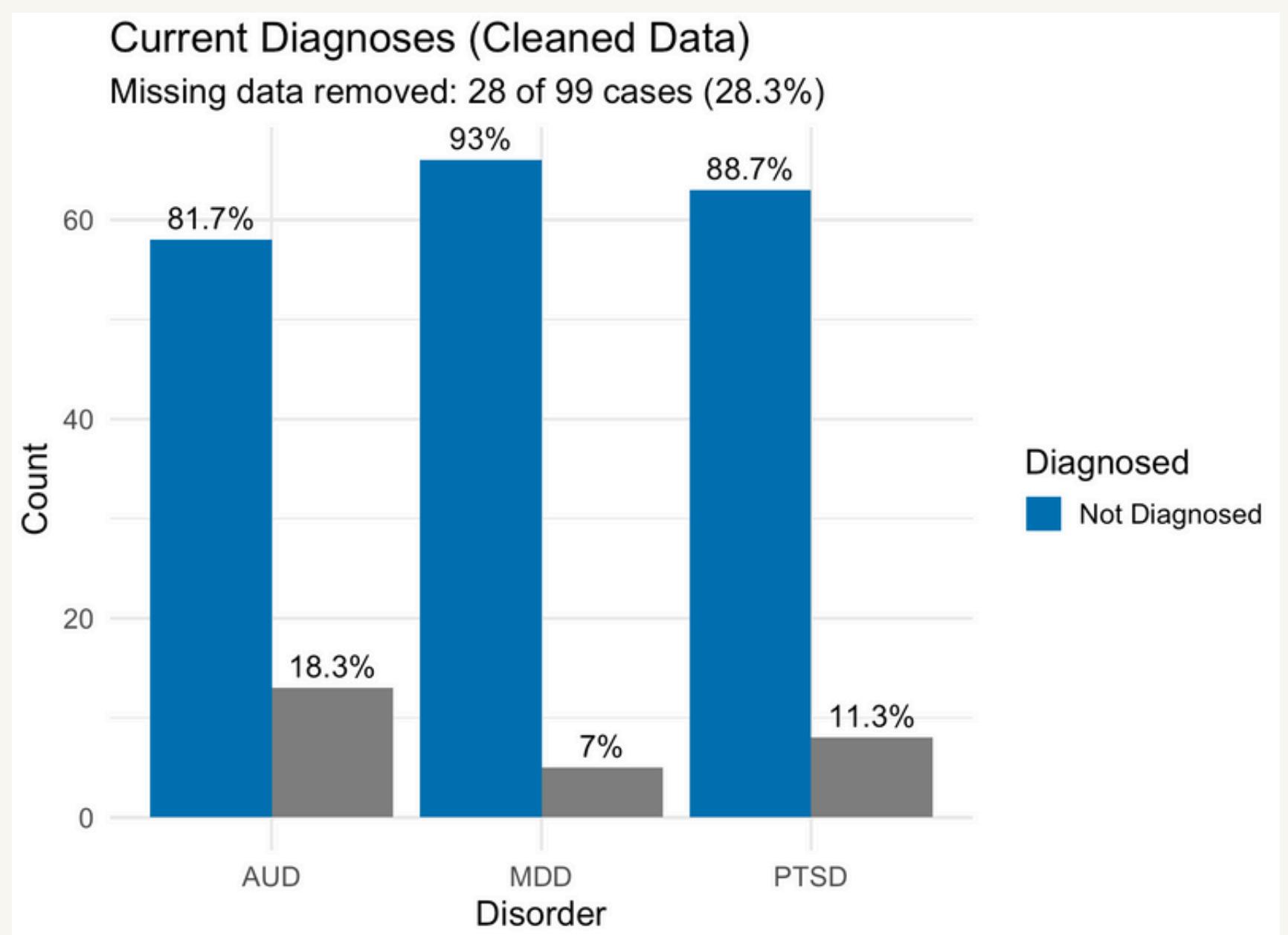


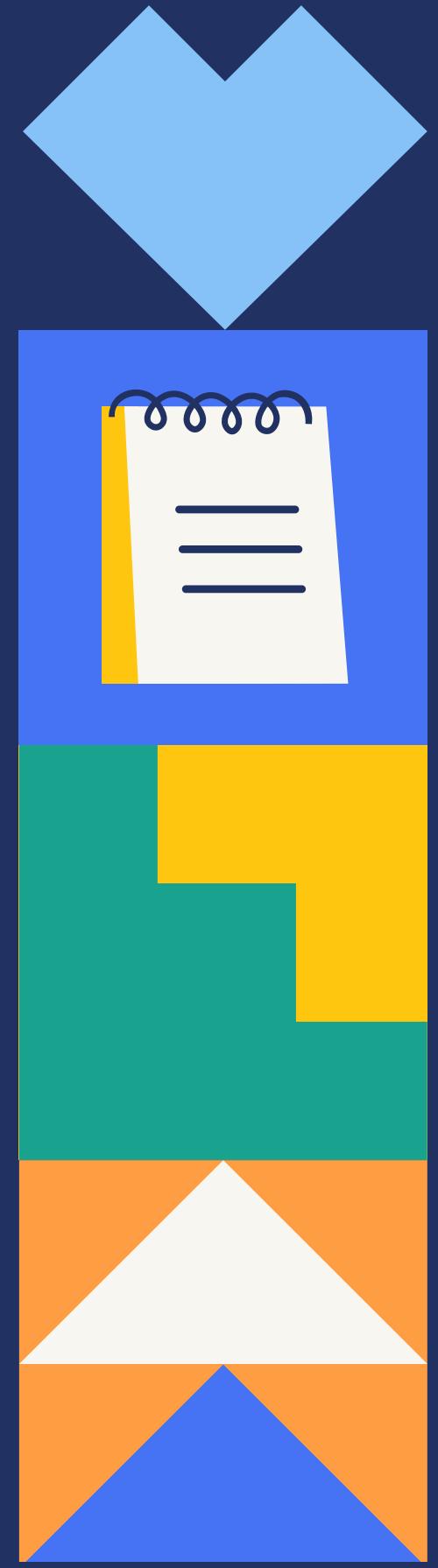
Add a subtitle to explain
missingness



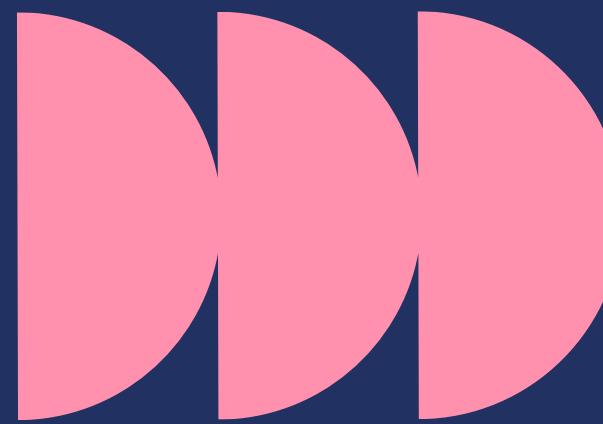
Presenting Frequencies

Graph of our diagnostic data with our clean dataset while being transparent about missingness





Outline



R Crash Course



Diagnosing Data Problems

- Visualising ✓
- Descriptives ✓
- Missing Data ✓
- Outliers



Curing Data Problems

- Handling missing data ✓
- Handling Outliers
- Visualising relationships

Diagnosing Problems: Continuous Data

Psychologists provided individuals with TBI a survey and then standardised their total score (we are neuropsychs so we love z scores).



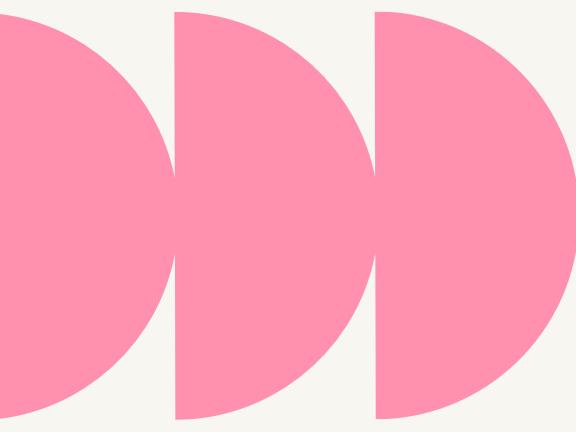
We will use this data to explore outliers now. This will be done visually.

Associations between transdiagnostic psychopathology dimensions and cognitive functioning

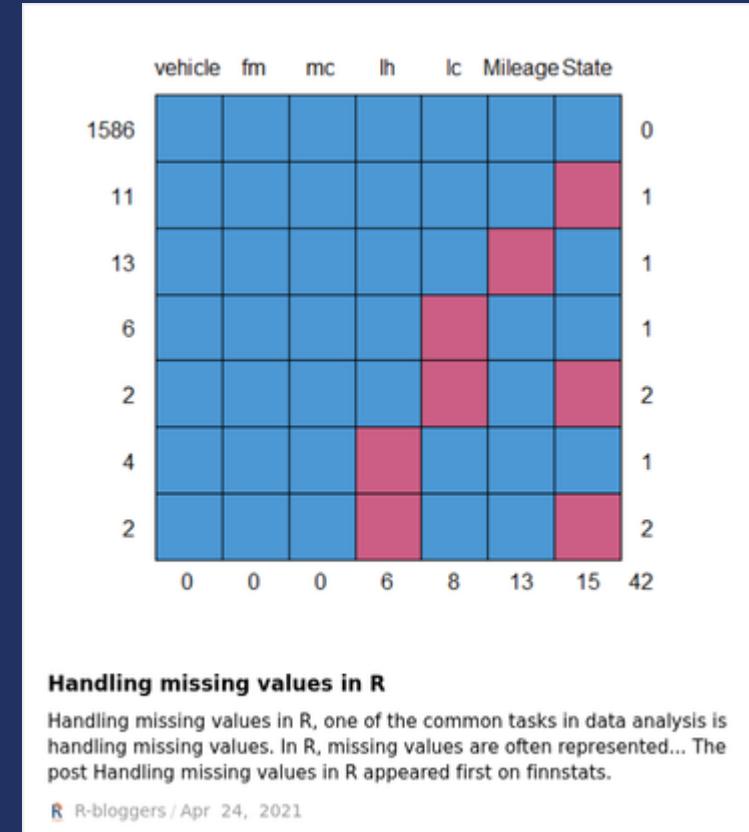
after traumatic brain injury: An application of the HiTOP-TBI model

Alexia Samiotis^{*1,2}, Jai Carmichael^{1,2}, Jao-Yue Carmintati^{1,2}, Amelia J Hicks¹, Jennie Ponsford^{1,2}, Kate Rachel Gould^{1,2}, Gershon Spitz^{1,2,3}

Continuous raw data



ID	z-score on psych scale
1	0.50
2	1.53
3	3.63
4	0.10

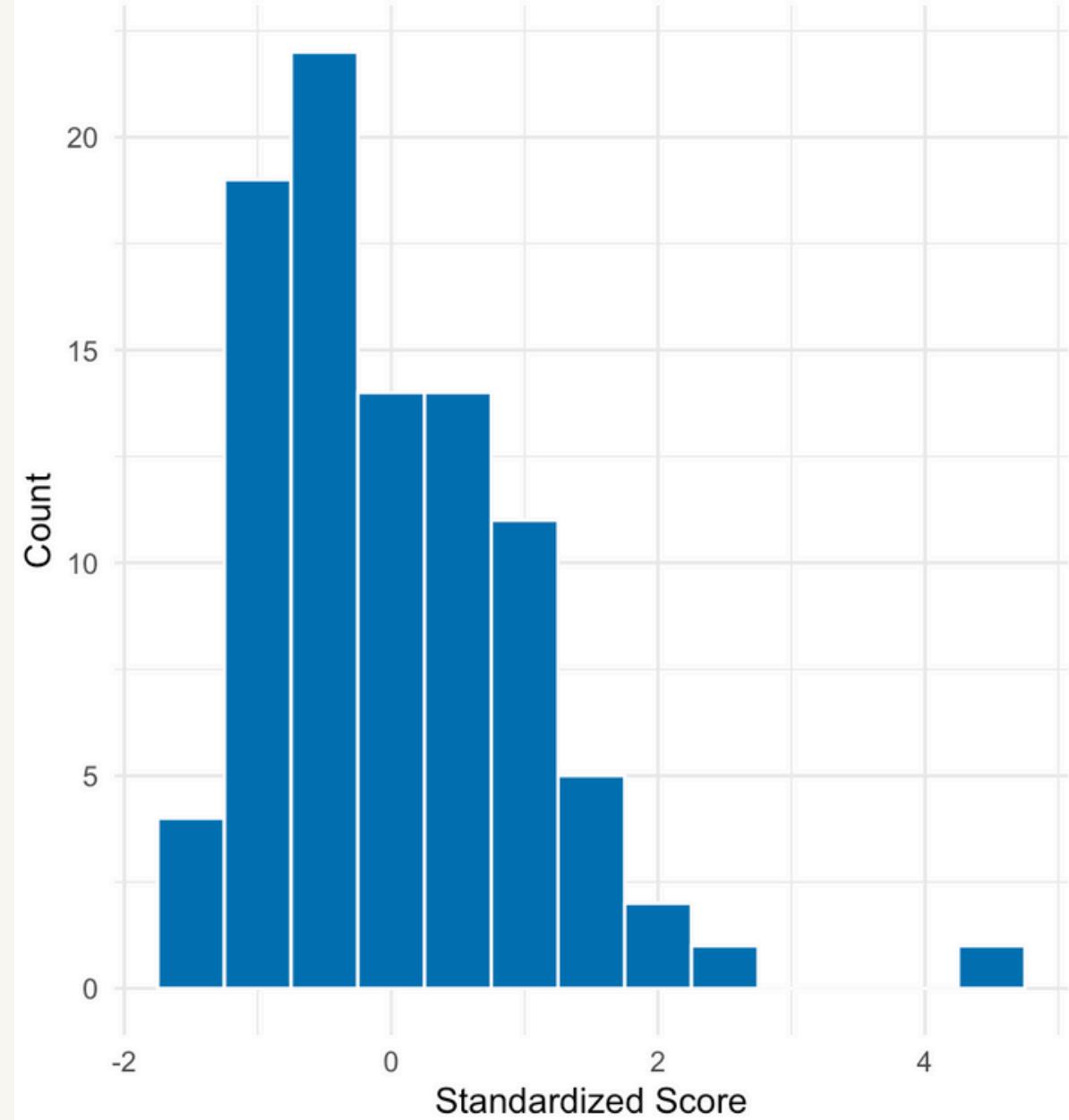


This is not the method I use to handle missing data, see other methods

Follow the steps outlined previously to import and clean the data by filtering NA values

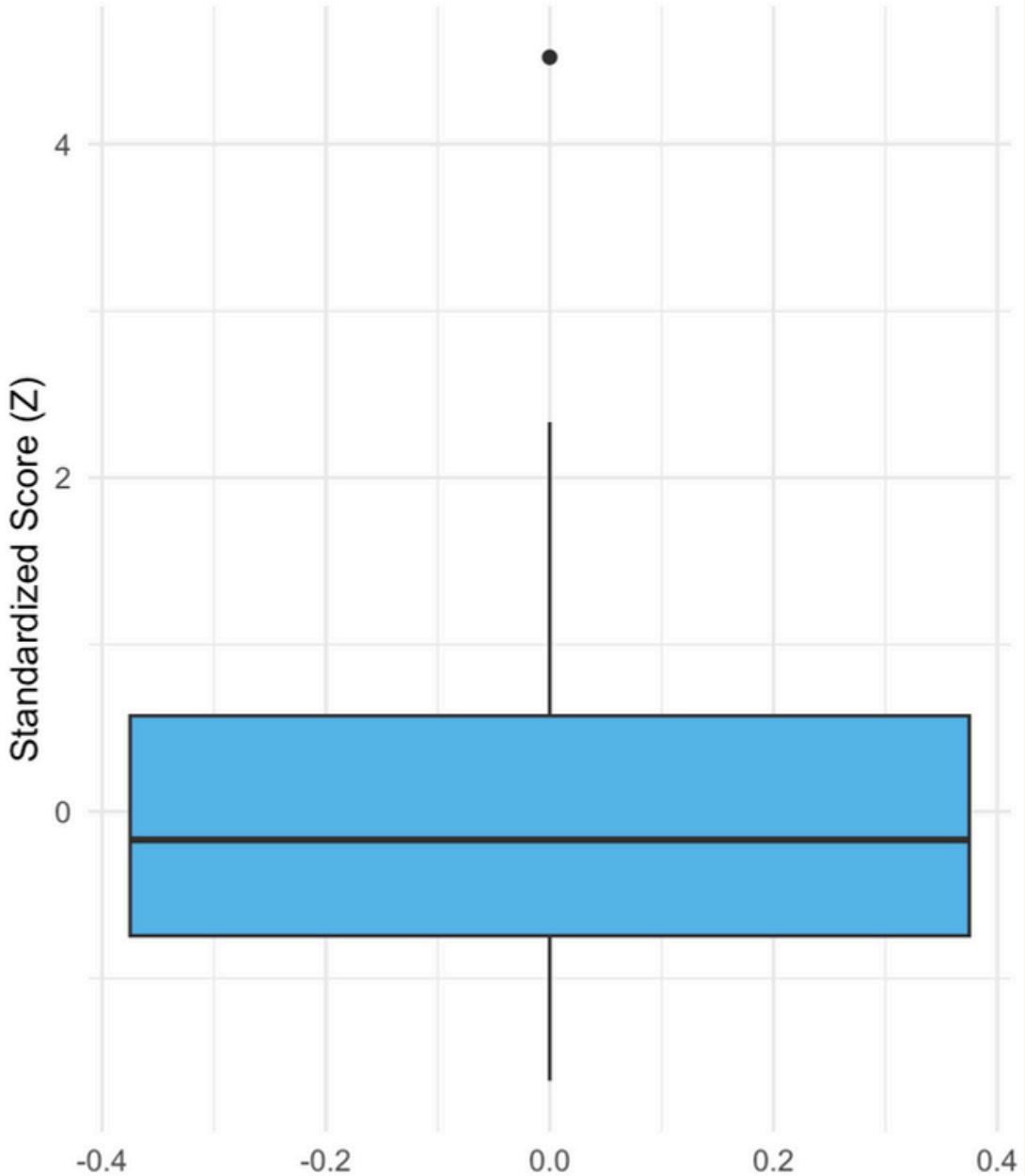
Visualisation for Outlier Identification

Distribution of General Problems



```
### Histogram – Distribution Overview
```{r}
ggplot(df, aes(x = General_Problems)) +
 geom_histogram(binwidth = 0.5, fill = "#0072B2", color = "white") +
 labs(title = "Distribution of General Problems",
 x = "Standardized Score",
 y = "Count") +
 theme_minimal()
```
```

Boxplot of General Problems

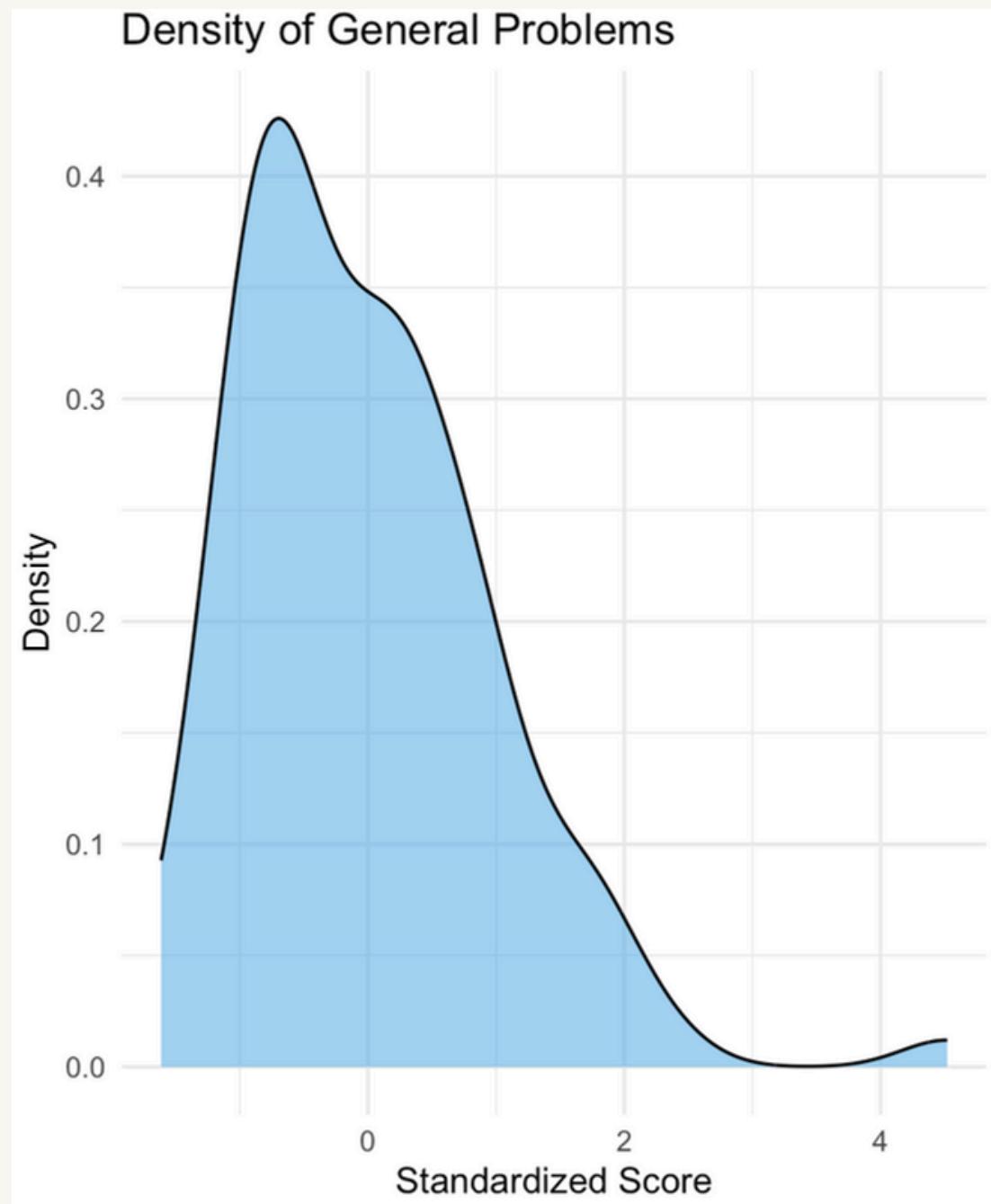


```
### Boxplot – Classic Outlier Detection
```{r}
ggplot(df, aes(y = General_Problems)) +
 geom_boxplot(fill = "#56B4E9") +
 labs(title = "Boxplot of General Problems",
 y = "Standardized Score (Z)") +
 theme_minimal()
```
```

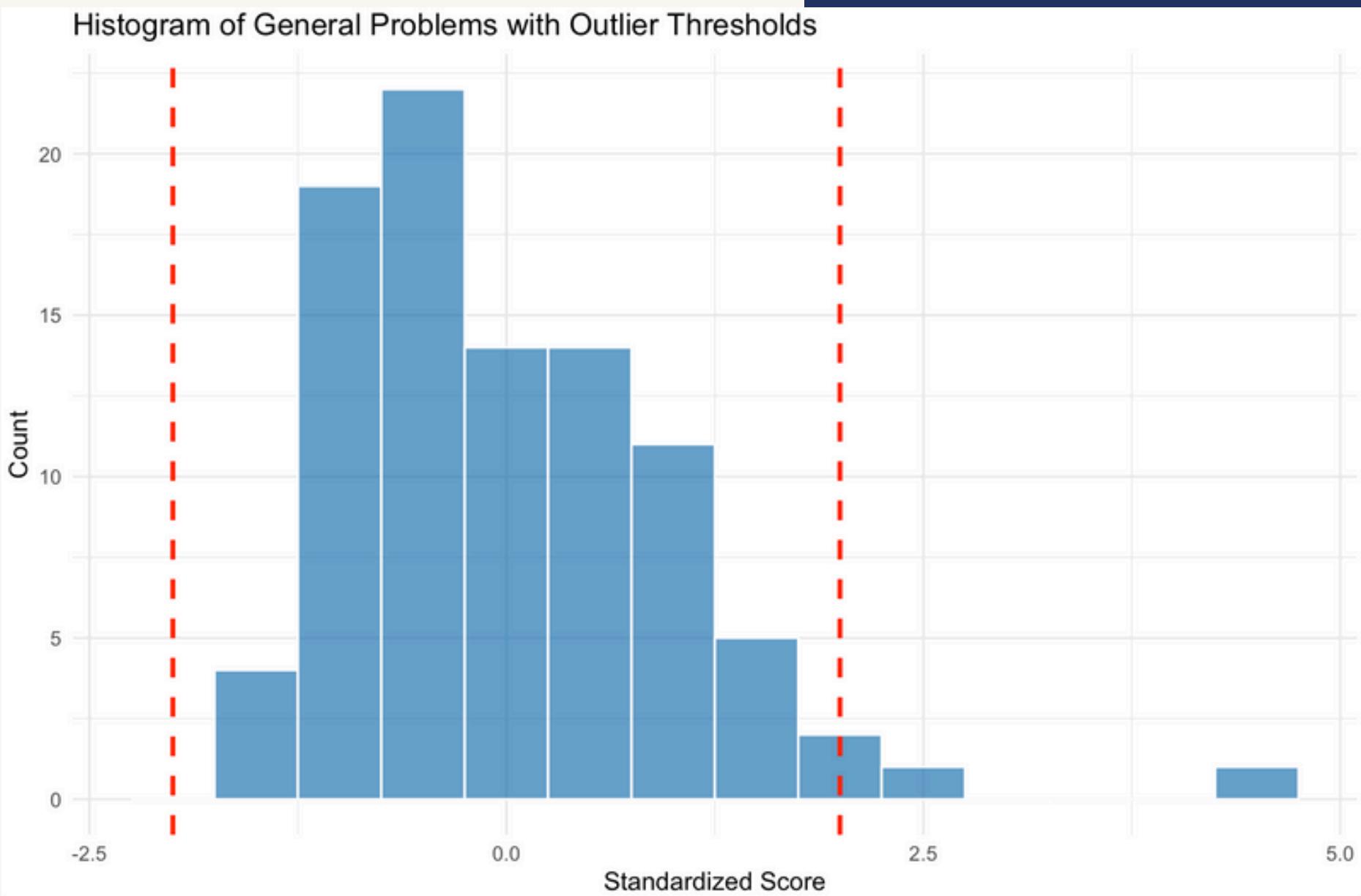
Values beyond ± 2 or ± 3 SD are often flagged as potential outliers

Visualisation for Outlier Identification

Values beyond ± 2 or ± 3 SD are often flagged as potential outliers



```
### Density Plot - Smoothed Distribution
```{r}
ggplot(df, aes(x = General_Problems)) +
 geom_density(fill = "#56B4E9", alpha = 0.6) +
 labs(title = "Density of General Problems",
 x = "Standardized Score",
 y = "Density") +
 theme_minimal()
```
```



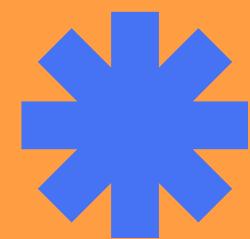
```
### Histogram or Density Plot with Threshold Lines
Shows where outliers fall relative to the main distribution.

```{r}
ggplot(df, aes(x = General_Problems)) +
 geom_histogram(binwidth = 0.5, fill = "#0072B2", color = "white", alpha = 0.7) +
 geom_vline(xintercept = c(-2, 2), linetype = "dashed", color = "red", size = 1) +
 labs(title = "Histogram of General Problems with Outlier Thresholds",
 x = "Standardized Score",
 y = "Count") +
 theme_minimal()
```
```

Two potential outliers were identified

Problem 2: Outliers

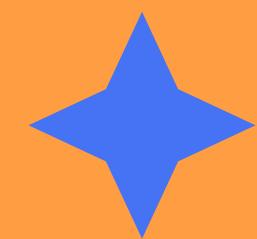
There are many ways to handle outliers once identified.



Outliers

Can be anything > 2SD from the mean

As identified on boxplots or other visuals

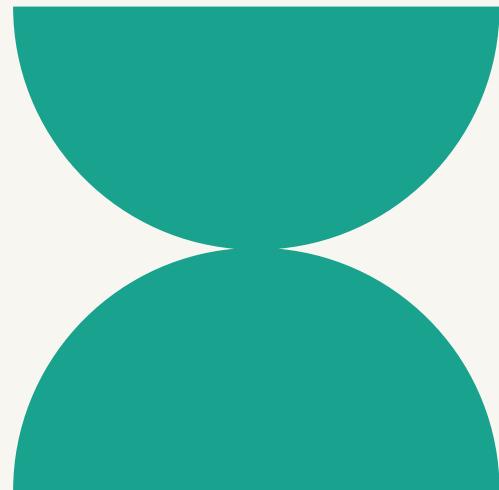


Solution

Leave them

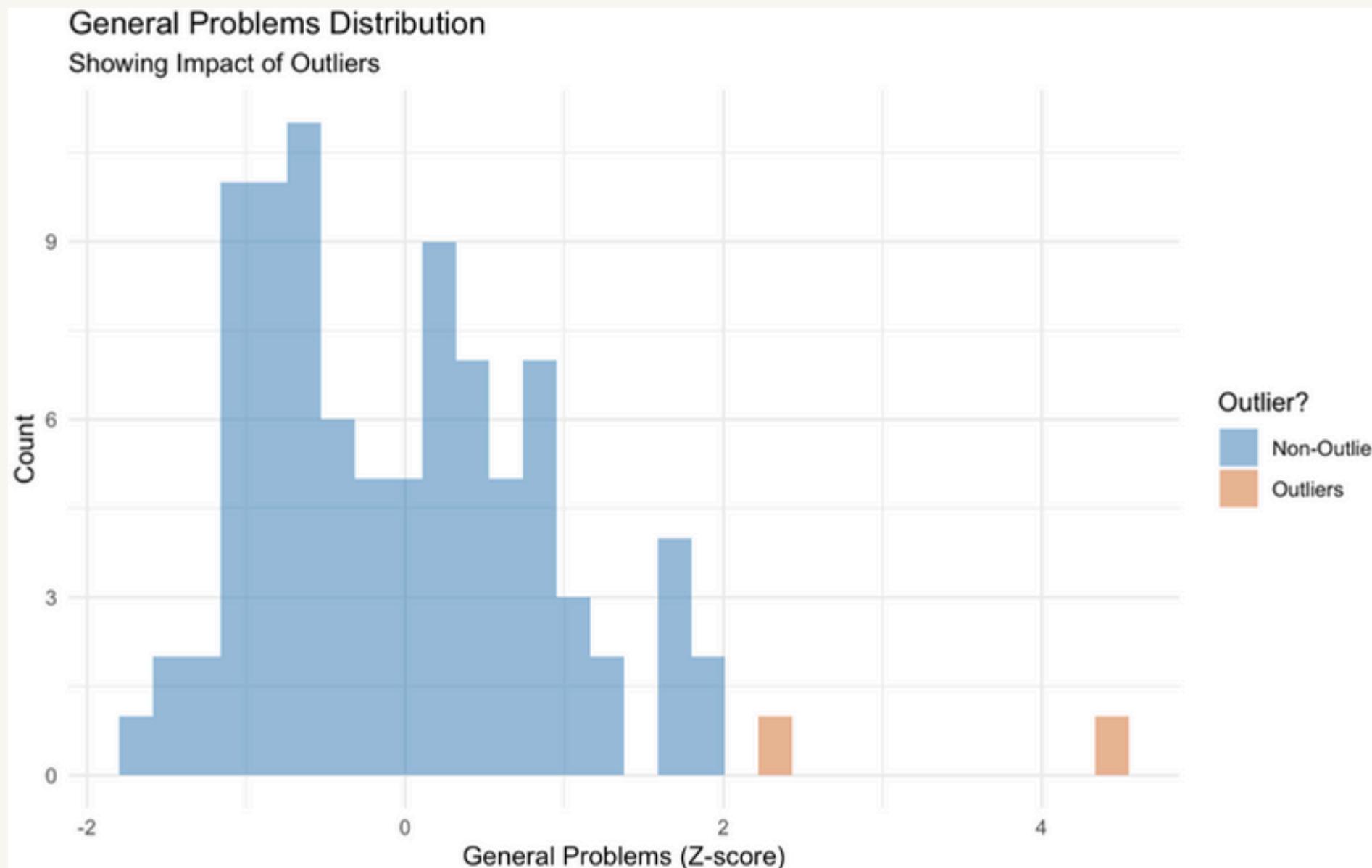
Remove

Winsorise (not covering)



Removing Outliers

We can compare the data with and without outliers before deciding what to do.



```
### Option 3: Run analyses with and without outliers and see if it affects it.  
```{r}  
df <- df %>%
 mutate(outlier_flag = abs(General_Problems) > 2)

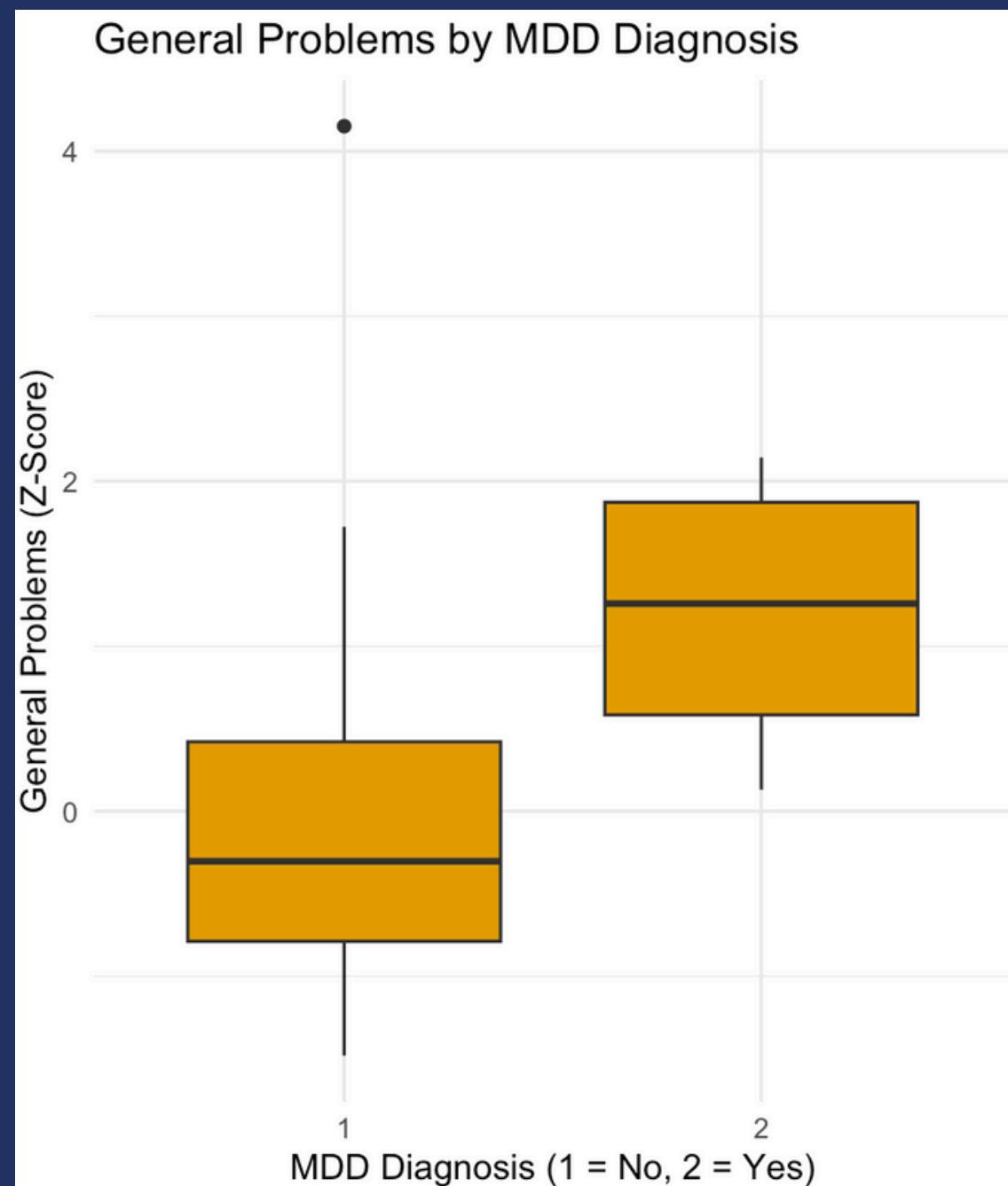
df %>%
 group_by(outlier_flag) %>%
 summarise(
 n = n(),
 mean_gp = mean(General_Problems, na.rm = TRUE),
 sd_gp = sd(General_Problems, na.rm = TRUE)
)

ggplot(df, aes(x = General_Problems, fill = outlier_flag)) +
 geom_histogram(position = "identity", alpha = 0.5, bins = 30) +
 scale_fill_manual(values = c("FALSE" = "#0072B2", "TRUE" = "#D55E00"),
 labels = c("Non-Outliers", "Outliers")) +
 labs(title = "General Problems Distribution",
 subtitle = "Showing Impact of Outliers",
 fill = "Outlier?",
 x = "General Problems (Z-score)",
 y = "Count") +
 theme_minimal()
```
```

Bonus: visualising relationships

We can do all the same techniques discussed while considering other important variables.

Interesting to present data this way as it shows how diagnostic labels don't neatly fit around all individuals with severe psychopathology after TBI.



This is consistent with some of the findings from my PhD - dimensional/continuous scales better capture psych problems than DSM diagnoses.

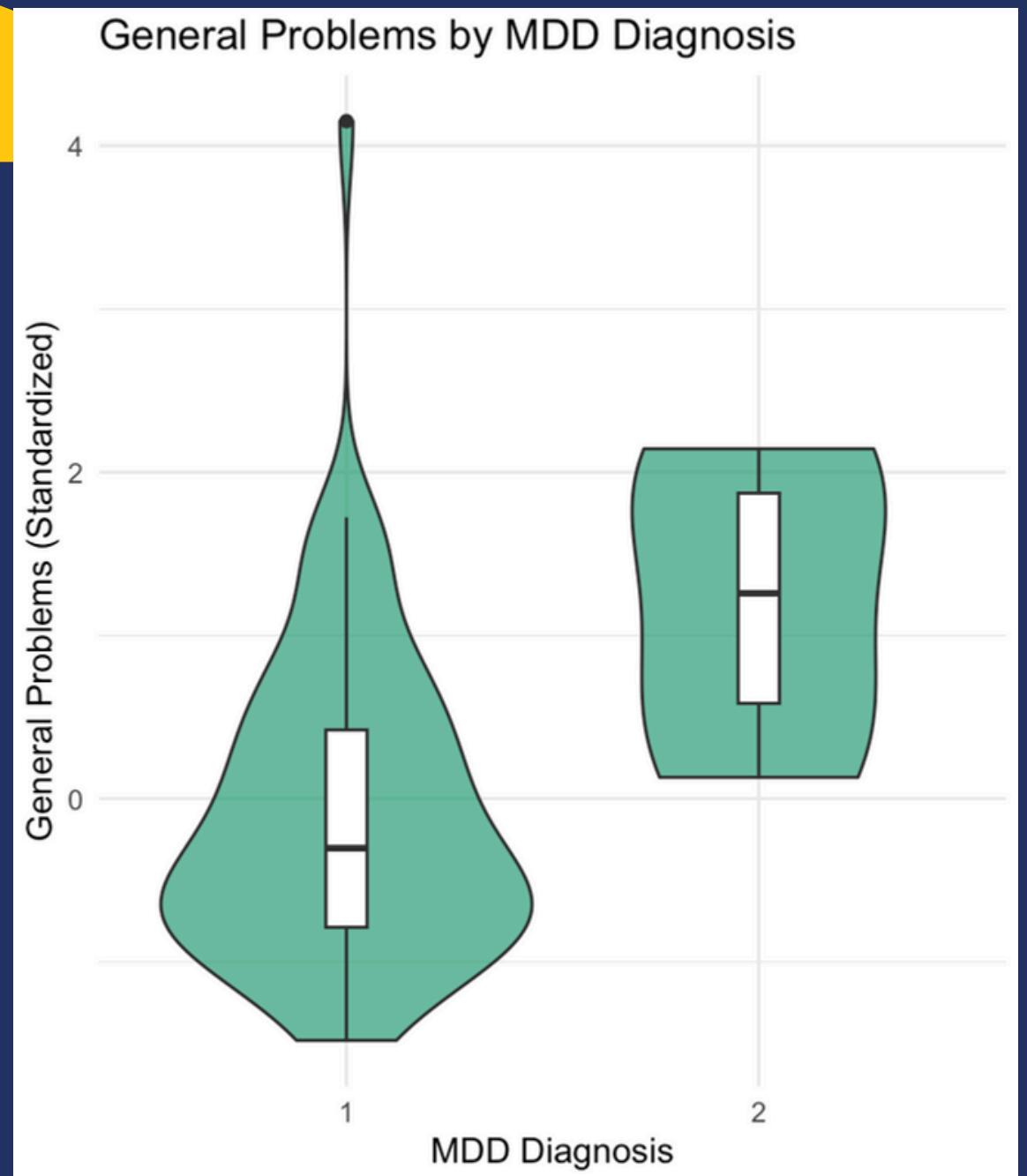
There is one outlier with an extreme (high psychopathology severity) score in the no diagnosis group

```
## Synergistic Visualisation: Using Two or More Variables
### Boxplots by Group – Compare Across Diagnoses or Gender
Example: By MDD diagnosis

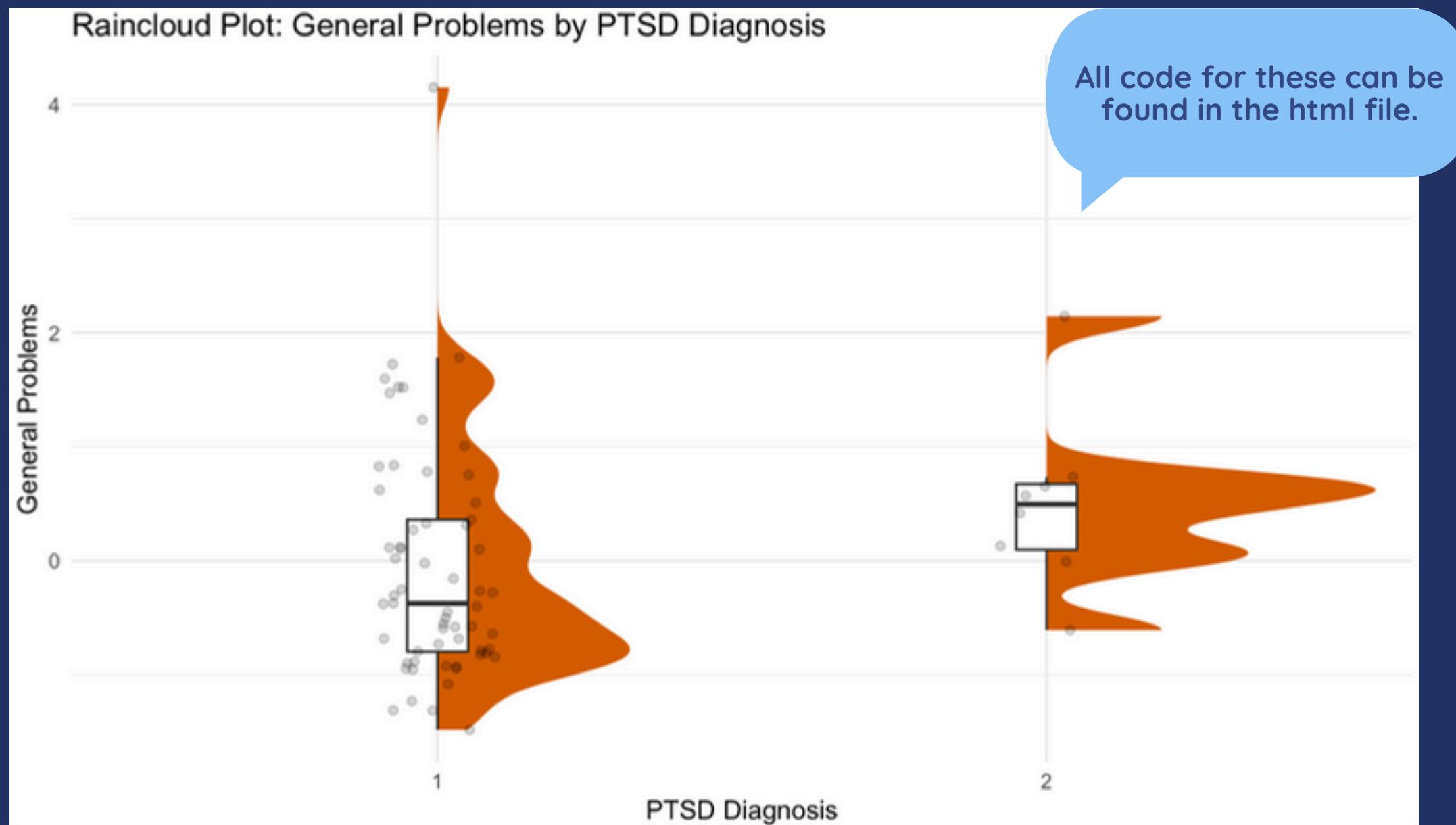
```{r}
df_clean$General_Problems <- scale(df_clean$a1_composite_score)
ggplot(df_clean, aes(x = factor(mini_mdd_current), y = General_Problems)) +
 geom_boxplot(fill = "#E69F00") +
 labs(title = "General Problems by MDD Diagnosis",
 x = "MDD Diagnosis (1 = No, 2 = Yes)",
 y = "General Problems (Z-Score)") +
 theme_minimal()
```
```

Bonus: visualising relationships

We can do all the same techniques discussed while considering other important variables.



```
### Violin Plot – Combine Density + Boxplot
```{r}
ggplot(df_clean, aes(x = factor(mini_mdd_current), y = General_Problems)) +
 geom_violin(fill = "#009E73", alpha = 0.7) +
 geom_boxplot(width = 0.1, fill = "white") +
 labs(title = "General Problems by MDD Diagnosis",
 x = "MDD Diagnosis",
 y = "General Problems (Standardized)") +
 theme_minimal()
```
```



```
### Raincloud Plot (Optional, Modern & Informative)
```{r}
library(ggdist)
ggplot(df_clean, aes(x = factor(mini_ptsd_current), y = General_Problems)) +
 stat_halfeye(adjust = 0.5, width = 0.6, .width = 0, fill = "#D55E00") +
 geom_boxplot(width = 0.1, outlier.shape = NA) +
 geom_jitter(width = 0.1, alpha = 0.2) +
 labs(title = "Raincloud Plot: General Problems by PTSD Diagnosis",
 x = "PTSD Diagnosis",
 y = "General Problems") +
 theme_minimal()
```
```

Thank you!

Associations between transdiagnostic psychopathology dimensions and cognitive functioning

after traumatic brain injury: An application of the HiTOP-TBI model

Alexia Samiotis^{*1,2}, Jai Carmichael^{1,2}, Jao-Yue Carmintati^{1,2}, Amelia J Hicks¹, Jennie

Ponsford^{1,2}, Kate Rachel Gould^{1,2}, Gershon Spitz^{1,2,3}

Supervisors

- Prof Jennie Ponsford
- Dr Gershon Spitz
- Dr Jai Carmichael

Questions or need any code?
alexia.samiotis1@monash.edu
or find me at the desks at MERRC
Thursdays/Fridays 😊

Link to code

