

Week 5: Survey

Charlotte Hadley

Topics for today

```
1 library(tidyverse)
2 library(gt)
3 library(readxl)
4 library(here)
5 library(janitor)
6 library(haven)
```

This week we're going to be discussing surveys
The goals for the lecture section of today

1. Identify what makes surveys

2. Correctly identify if data is "long" or "wide"

3. Understand how to use the {tidyr} pivot functions

4. Understand how to use the tidyverse for cleaning survey data

We'll likely continue some of the lecture notes

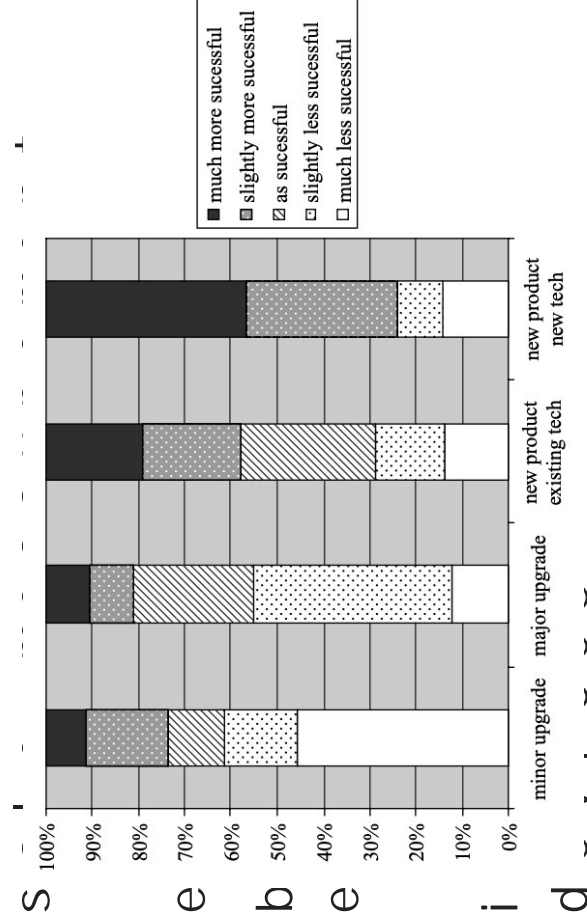
Surveys are an indispensable healthcare

There are lots of absolute measures in healthcare [treatment outcomes] ...

- Patient wait times • ... for understanding patient
- Morbidity • ... for tracking patient
- Biological samples • ... for medical trials
- Physiological health • measuring medical
- Device-based measurements
- Anthropometric measurements to understand the patient
- Sensory measurements

Surveys are an indicator healthcare

Surveys might be the
we take in a study
NHS patient experience
"The importance of urban
and urban ecosystem
COVID-19 pandemic"
How successful do medical
companies rate their



Source: Eatock, et al

Surveys are an indispensable healthcare

Surveys might be the *single* most effective way to take a study context for other measures

NHS patient experience *take* surveys

"The importance of *urban* health studies that take and urban ecosystem *assessment* but the for COVID-19 *pandemic*"

How successful do medical *technology* measures companies rate *their* *de* *survey* *logical* *survey* *survey*

Designing effective

We don't have enough time to go to deep surveys
probably an entire undergraduate course in

So why are we looking at surveys?

Designing effective

There are ~~from~~ ~~the~~ ~~tips~~ ~~in~~ However these tips do n
designing surveys I waerft ctt d vœœsurvey

In effective survey:

These "tips" are geared towards a design g
surveys where you can easily analyse the
data after running the survey. ~~is asking~~.

- Questions are unbiased leading.

The best way to test
your survey³ is pretest

Designing effective

Getting survey data isn't the only task a manager requires after a week of fieldwork and tasks/skills. Wickham in 2014.

We will use survey data as an introduction to wrangling:

- Pivoting data between wide and long formats
- Joining datasets [in the same way that SQL
- Wrangling survey questions with multiple c
- Wrangling survey questions that capture mu

Quantitative vs
Qualitative measurements?

"Closed question" vs
which is which?

Surveys can capture

Quantitative measures of internet use
collected with closed-ended questions

- Do you ^{filed} use tracking device
 - Yes / No
- How often do you wear what you ^{would} like to recommend telling
 - Every day / Some days / A few days / Rarely
- Since owning a tracking device do you feel like you know more about your activity?
 - Strong agree / Agree / Neither agree nor disagree / Disagree / Strong disagree

Surveys can capture data

Ideally we would always use closed-form questions but unfortunately this isn't always the case. If Survey

Question" Where do you get most of your information from?"

Question information text

Example responses:

What do you think about

Q. Do you support the idea that charities should not pay tax?

- ☐ Yes
- ☐ No
- ☐ Don't know

This is a question for
bite-size guides that we write
questionnaires

Survey Mode

Methods of survey data

There are different models considered for survey data collection: a cfoel d teicntg osnu: ravfe⁷

- Online (open/close)
 - Telephone
 - Mail
 - Face-to-face
 - Paper (observed)
 - Mixed-mode
 - Same survey and modes
 - Multi-phase survey modes
-
- by mode⁸
- Respondents answer questions by mode⁸
- Respondent demographic and survey topic
-
- Inaccurate state-level 2016 US elections are have been strongly biased over-representation of certain groups

Survey Size

Survey Size

Many survey tools provide population size

Population	Required sample size
528 - people who have been to space	223
10,490 - athletes at the London 2012 Olympic Games	371
110,000 - wine growers in France	383
5,300,000 - all Hebrew speakers	384
50,000,000 - everyone who's bought Michael Jackson's "Thriller"	384
1,344,130,000 - everyone in China	384

- These are reasonable population size
- studies" but for other variables using a correlational design. Results showed that sex (i.e., male), low minority exploration, and low other-focused orientation were risk factors for primary psychopathology in terms of the importance of emerging adulthood development.
- For pre-test surveys
- to use sample at least
- per Bernger et al

Source: 12

LIKERT Scales (I)

In a LIKERT scale responses are given a score.

response	
Strong c	1
Disagree	2
Neutral	3
Agree	4
Strong a	5

It is absolutely meaningless of Ta sldaklEERb ut peo do it.

When you create your credit in a

question	response
Feel fidæmt in	Strong a
Feel fidæmt in	Neutral
Feel fidæmt in	Disagree
Feel fidæmt in	Strong d

LIKERT Scales (LI)

response	
Strong c	1
Disagree	2
Neutral	3
Agree	4
Strong a	5

question	response
Feel fidæmt in ·Strong a	
Feel fidæmt in ·Neutral	
Feel fidæmt in ·Disagree	
Feel fidæmt in ·Strong d	

Our original questionnaire had a variable that required responses.

```
1 cat("Strong disagree < Disagree ... < Strong agree")
```

Strong disagree < Disagree ... < Strong agree

When we convert the questionnaire scores are the interval data

But is this accurate
Is the difference between "Strong disagree" and "Disagree" the same as that between "Neutral Agree"?

LIKERT Scales (III)

response	
Strong c	1
Disagree	2
Neutral	3
Agree	4
Strong a	5

If you want to compare multiple sliders then the meaning to the "media

You could perform fact is well explained¹⁴by

question	response
Feel fidæmt in ·Strong a	
Feel fidæmt in ·Neutral	
Feel fidæmt in ·Disagree	
Feel fidæmt in ·Strong d	

When designing your s also directly ask res instead of trying to + "neutral" means.

Overall, do you think using the tidyverse?

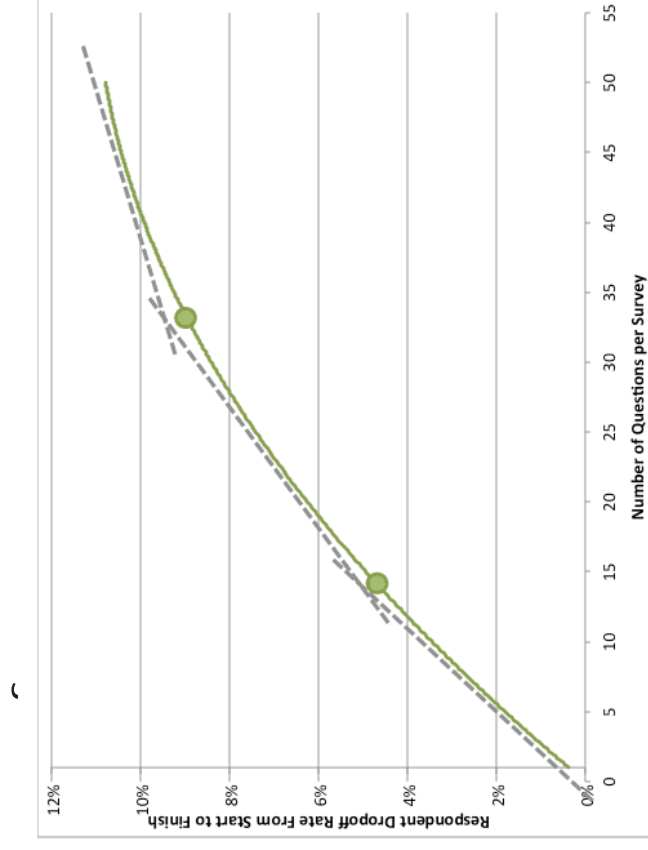
Missing data in

Missing data in surveys

When designing a survey we ideally want responses

However, research¹⁵ has found that response rate (or survey abandonment) correlated with survey length. Although, the effects by

- survey mode
- survey reward
- are questions skippable



Source: SurveyMonkey¹⁵

Missing data in surveys

There are three kinds of missing data distribution

- Missing Not at Random (MNAR) missingness is related to the response
- Missing at Random (MAR) missingness of response is related to question characteristics, but not to the response itself
- Missing Completely at Random (MCAR) missingness is unrelated to any characteristics of the respondent (e.g. extra-marital relationships)

Data st MANTRA indicates study design.

However, of the three, it is indeed MCAR that is the most difficult to detect in your data.

Missing data in surveys

There are three kinds of missing data distribution

- Missing Not at Random (MNAR) is not rare
- Missing at Random (MAR) would call this "conditional" on observed data
- Missing Completely at Random (MCAR) by including information by order
 - Missingness might be different in demographic males are less likely to depression surveys.

Data stack provides useful methodology for imputing missing data

Missing data in surveys

There are three kinds of missing data distribution

- Missing Not at Random
 - Missing at Random
 - Missing Completely at Random
- MCAR (c_{MCAR}) is harmonic with fluctuation study for randomness.

- Randomly sampling across questions for each participant

Planned missingness depends on fairly advanced statistical methodologies.

Survey tools

Survey tools

There's a plethora of survey tools available.

Many of these tools provide free tiers but [and for GDPR compliance].

One thing I don't like about these tools is that they don't have a good format - if you're using them for a long time, it's a pain to format.

Let's look at the different types of survey tools.

Tidying multiple choice questions with R

This survey is being duplicated in Google Forms, Survey Monkey and Qualtrics. We are then writing a blogpost at rfortherestofus.com/blog about how to tidy multiple choice question datasets with R code.

[Sign in to Google](#) to save your progress. [Learn more](#)

Select all the things you've done in the past 24hours.

☐ Slept

☐ Eaten food

☐ Cooked food

☐ Gone to work

☐ Commuted for work

☐ Relaxed with a hobby (TELL US THE HOBBY BY TYPING IN THE OTHER FIELD)

☐ Other:

This survey is being duplicated in Google Forms, Survey Monkey and Qualtrics. We are then writing a blogpost at rfortherestofus.com/blog about how to tidy multiple choice question datasets with R code.

* 1. Select all the things you've done in the past 24hours?

☐ - Slept

☐ - Eaten food

☐ - Cooked food

☐ - Gone to work

☐ - Commuted for work

☐ - Relaxed with a hobby (TELL US THE HOBBY BY TYPING IN THE OTHER FIELD)

☐ Other (WHAT'S YOUR HOBBY?)

Tidying multiple choice questions with R
This survey is being duplicated in Google Forms, Survey Monkey and Qualtrics. We are then writing a blogpost at rfortherestofus.com/blog about how to tidy multiple choice question datasets with R code. the question text
Select all the things you've done in the past 24hours.

Slept

Eaten food

Cooked food

Gone to work

Commuted for work

Relaxed with a hobby (TELL US THE HOBBY BY TYPING IN THE OTHER FIELD)

Other (WHAT'S YOUR HOBBY?)

Survey tools / software

Because there's a huge variety in data representation, train your general purpose data wrangling skills

In the lecture we'll look at several real-world wrangling.

In the workshop you'll look at simpler data learn now

Survival for mats

- Excel files: Almost all files we have seen so far are arranged in a way that the data is arranged in a way that is easy to read. We use the {readxl} package for reading Excel files. If data is encoded with cell colour you'll need the readxl package.
- .csv files: Most tools provide a way to export data as a .csv file. They're likely to be well formatted. Files like .csv are read into R with the read.csv() function.

Survey datasets

Survey datasets (I)

We'll be looking at datasets during this week:

- Emerging Adulthood Measured ¹⁷at Multiple In-
 - This is the 2nd instance of a large scale
 - Learn about the 1st instance (and how ¹⁸ the
 - The actual survey question ¹⁹ is [here](#) ²⁰ available
 - The survey dataset ²¹ is available ²² for download
 - The actual survey ²³ data ²⁴ is available ²⁵ for download via a URL.

Survey datasets (II)

We'll be looking at datasets during this week:

- Emerging Adulthood Measured ¹⁷
- Public Attitudes to Commercial Access to Health Data
 - The WI comes from the Ipsos MORI to see access to health data
 - **This study can be read here**
 - The survey questions can be found at the
 - The dataset is openly available ²⁰ but it has by account.

Survey datasets (II)

We'll be looking at datasets during this week:

- Emerging Adulthood Measured ¹⁷
- Public Attitudes to Commerce ²⁰
- British Election Study 2019
 - Since 1964 a post-election survey has been motivations and the impact of political
 - Data for all surveys is available at data.britishteasectiional-data/
 - The 2019 questionnaire and information on britishel.eccotmi/odnasttau-doyb ject / 2019-british-elrandom-probability-survey/ ²²
 - The 2019 survey data is available ⁵ but must be added to your UK Data Service account.

Survey datasets (IV)

We'll be looking at 3 datasets during this week:

- Emerging Adulthood Measured at Multiple Intervals
- Public Attitudes to Commercial Access to Human Tissues
- British Election Study 2019



Task: Setup our proj

SLIDE 1 OF 3

1. Create a new `epnrgo7j2e1c8t-waeakl-e5d_surveys`
2. Add a subfolder `data` to the dataset.
3. Create a `README` for each of the studies
 - `emerging-adulthood.Rmd`
 - `commercial-access-to-health-data.Rmd`
 - `british-election-study-2019.Rmd`



Task: Obtain Emerging data

SLIDE 2 OF 3

1. Open the emerging-films website. Rmd

There are two prompts a / t / o'3 ft .h i a d / qwe q m k e d :

- The codebook
- The dataset

2. Add a code chunk to the Rmd file.

3. Add this code chunk to the Rmd file.



Task: Read in Emergency data

SLIDE 3 OF 3

When we read in datasets we should always raise the ~~file~~ with object names that indicate this

```
1 adult_hood_raw_data <- read_excel("data/emerging-adulthood_data.xlsx")  
2 adult_hood_raw_codebook <- read_excel("data/emerging-adulthood_codebook.xlsx")
```


Messy column names (

Most find a\$ you' ll work with will have messy c
with:

```
1 glimpse(adulthood_raw_codebook)
```

```

Rows: 328
Columns: 6
$`Variable Name`      <chr> "StartDate", "EndDate",
"Status", "Progress", "Du...
$`Question text`      <chr> "Start Date", "End Date",
"Response Type", "Progr...
$...3                 <chr> "n/a", "n/a", "n/a", "n/a",
"n/a", "n/a", "n/a", ...
$ responses           <chr> "qualtrics variable",
"qualtrics variable", "qual...
$...5                 <dbl> NA, NA, NA, NA, NA, NA, NA, NA,
NA, NA, NA, NA, NA, N...
$`Survey Question ID` <chr> "
`\ImportId\":"startDate\\"", "{\ImportId\":"..."

```

```
1 adult_hood_raw_codebook %>%
2   select(`Question text`)
```

```
# A tibble: 328 x 1
  `Question text`
<chr>
1 Start Date
2 End Date
3 Response Type
4 Progress
5 Duration (in seconds)
6 Finished
7 Recorded Date
8 Response ID
9 Recipient Last Name
10 Recipient First Name
# ... with 318 more rows
```

The easiest way to close a window is to click the close button in the top right corner of the window.

Messy column names (

But understand how these two datasets relate

There are `names_raw_codebook` links the `names_raw_data`

If we can update the `names_raw_codebook` with no longer this instantly clean the codebook column names

```
1 adult_hood_raw_codebook <- read_excel("data/emerging-adult_hood_codebook.xlsx") %>%  
2 clean_names()
```

I DEA Questions (I)

Let's take a look at the *thesour* data *ways* et *questions*

	Strongly disagree	Somewhat disagree	Somewhat agree	Strongly agree
Is this period of your life a time of defining yourself?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Is this period of your life a time of deciding your own beliefs and values?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Is this period of your life a time of high pressure?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Is this period of your life a time of many possibilities?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Is this period of your life a time of gradually becoming an adult?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Is this period of your life a time of exploration?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Is this period of your life a time of feeling stressed out?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Is this period of your life a time of feeling adult in some ways but not others?	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

```
1 adulthood_codebook %>%
2 filter(str_detect(question_text, "defining yourself"))
```

```
# A tibble: 1 x 6
  variable_name question_text x3 response x5 survey_id
<chr>          <chr>          <chr> <chr> <dbl> <chr>
1 IDEA_5      Is this period of your life a time ... This... 1-stro... NA "{\\"Im..
# ...with abbreviated variable names `responses`, `survey_question_id`
```

```
1 idea_responses_raw <- adulthood_raw_data %>%
2 select(ResponseId, starts_with("IDEA_"))
3 idea_responses_raw
```

```
# A tibble: 3,182 x 9
  ResponseId IDEA_1 IDEA_2 IDEA_3 IDEA_4 IDEA_5 IDEA_6 IDEA_7 IDEA_8
<chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 R_BJN3bOqi 1zUM d3      3      4      4      3      4      4      4
2 R_2TGbi BXm txyw sD      4      4      4      3      4      4      4
3 R_12G7bl qN2wB2N65      4      4      4      4      4      3      3
4 R_39pl dNoon8CePfP      4      4      3      3      4      4      4      4
5 R_1Q Kb2LdJo1Bhvv      4      4      3      3      4      3      3      4
```

6	R_pmwDTZyCyCycXwB	3	4	3	3	4	4	3	2
7	R_2QjhOh3wxTj Zj KP	4	3	4	4	3	4	3	3
8	R_2CfdmFwINTI i v4e	4	3	3	4	3	2	2	3
9	R_24kJ PxVOxM\$hN3Q	4	4	3	4	4	4	3	3
10	R_3fvOvHsw6AvJPK	4	4	2	4	4	4	3	4

... with 3,172 more rows

I D E A Q u e s t i o n s (I I)

Can you tell me what dataset?

```
1 idea_responses_raw %>%  
2 head() %>%  
3 gt()
```

```
ResponseId  I D E A I D E A I D E A I D E A I D E A  
R_BJN3bQqi1; 3 4 4 4 3 3 4  
The adulthood_raw_codes book tells us that  
" 4" encodes " Strong agree"  
R_12G7bIqN2\ 4 4 4 4 4 4  
R_39pldNoon{ 4 4 4 3 3 4  
We now need a way to find out what  
these columns at once contain you suggest  
R_pmwwTCyCy 4 4 4 4 4 4
```

There are two methods I can think of:

- One method we've already used in coding
- One method we'll be introducing today

I DEA Questions (I I I)

- Using gross target multiple columns at once

```
1 idea_responses_raw %>%
2   mutate(across(starts_with("IDEA_"),
3     ~case_when(.x == 1 ~ "Strong disagree",
4       .x == 2 ~ "Disagree",
5       .x == 3 ~ "Agree",
6       .x == 4 ~ "Strong agree")))
```

- Using pivot_longer() transform this from wide to

```
1 idea_responses_raw %>%
2   pivot_longer(starts_with("IDEA_")) %>%
3   mutate(value = case_when(value == 1 ~ "Strong disagree",
4     value == 2 ~ "Disagree",
5     value == 3 ~ "Agree",
6     value == 4 ~ "Strong agree"))
```

Using pivot_longer() the fitted dataframe
{ ggplot2 }

Wide vs Long data (1)

In wide datasets variables are listed in each row in a unique column. observation.

PersonWeight			
Bob	32	16	18
Alice	24	15	17
Steve	64	14	16

PersonAge	
Bob	Age 32
Bob	Weight 16
Bob	Height 18
Alice	Age 24

However, datasets might be For instance, year is interpreted as years for multiple columns. data.

countryyear	
UK	Supermar2002023
UK	Shopping4C4246
US	Supermar303638
US	Shopping8C9C98

countryyear	
UK	Supermar20020
UK	Supermar20020
UK	Supermar20023
UK	Shopping2004C
UK	Shopping20042
countryyear	

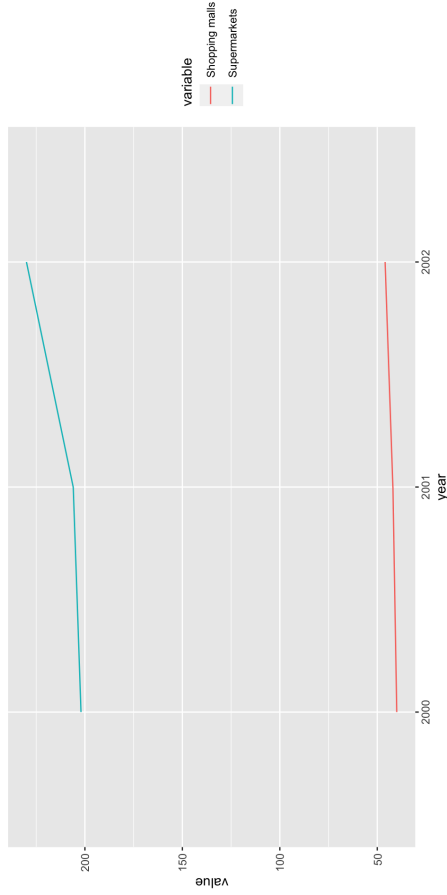
UK	Shoppin	200	4
US	Supermar	200	30
US	Supermar	200	36
US	Supermar	200	38
US	Shoppin	200	80
US	Shoppin	200	90
US	Shoppin	200	98

Wide vs Long data (I)

The `ggplot2` package requires data

country		year	value
UK	Supermarket	2000	20
UK	Supermarket	2001	20
UK	Supermarket	2002	23
UK	Shopping	2000	40
UK	Shopping	2001	42
UK	Shopping	2002	46
US	Supermarket	2000	30
US	Supermarket	2001	36
US	Supermarket	2002	38
US	Shopping	2000	80
US	Shopping	2001	90
US	Shopping	2002	98

```
1 long_shops_data %>%  
2   filter(country == "UK") %>%  
3   ggplot(aes(x = year,  
4             y = value,  
5             group = variable,  
6             color = variable)) +  
7   geom_line()
```



pivot_wider() and pivot_longer()

The `pivot_wider()` and `pivot_longer()` functions are for transforming format and wide format, respectively

```
1 tribble(
2   ~Person, ~Age, ~Weight, ~Height,
3   "Bob",    32L,    168L,    180L,
4   "Alice",  24L,    150L,    175L,
5   "Steve",  64L,    144L,    165L
6 ) %>%
7 pivot_longer(cols = Age:Height) %>%
8 gt()
```

Person	Age	Weight	Height
Bob	32	168	180
Bob	24	150	175
Bob	64	144	165
Alice	24	150	175
Alice	64	144	165
Steve	64	144	165

Note that `away` of `at` `days` `select` it `on` `target` `columns`

Pre 2020 `ts` `per` `read` `data` `get` `her` (The `se` `func` `t` `it` `and` `sy` `ra` `r` but `are` `con` `si` `der` `e` `pi` `s` `o` `p` `e` `r` `f` `cu` `en` `d` `ce` `ti` `ob` `ns` `.` `the`

I DEA Questions (I V)

We can now transform our actual dataset into

```
1 i_dea_responses_raw %>%  
2   pivot_longer(starts_with("I DEA_"))
```

The remaining step is to use the case_when()

```
1 i_dea_responses_long <- i_dea_responses_raw %>%  
2   pivot_longer(starts_with("I DEA_")) %>%  
3   mutate(value = case_when(value == 1 ~ "Strong disagree",  
4                             value == 2 ~ "Disagree",  
5                             value == 3 ~ "Agree",  
6                             value == 4 ~ "Strong agree"))  
7  
8 i_dea_responses_long %>%  
9   head() %>%  
10  gt()
```

ResponseID	name	value
R_BJN3bQqi	I DEA Agree	
R_BJN3bQqi	I DEA Strong	
R_BJN3bQqi	I DEA Strong	
R_BJN3bQqi	I DEA Agree	
R_BJN3bQqi	I DEA Strong	
R_BJN3bQqi	I DEA Strong	

IDEA Questions (V)

We need to match up the actual question

To achieve this we're mutating join

It's worth while mentioning this is a skill you would use in SQL.

If you're comfortable doing this then you'll be comfortable with basic SQL.

Source:

<https://github.com/kadenbuie>

I D E A Q u e s t i o n s (V I)

```

1 idea_responses_long %>%
2 head() %>%
3 gt()

```

ResponseId name value

R_BJN3bQqi IDEA Agree

R_BJN3bQqi IDEA Strong

R_BJN3bQqi IDEA Strong

R_BJN3bQqi IDEA Agree

R_BJN3bQqi IDEA Strong

R_BJN3bQqi IDEA Strong

Let's extract the variable labels from the codebook. Note that the column data is different

```

1 idea_question_labels <- adulthood_raw_codebook %>%
2   filter(str_detect(variable_name, "IDEA_")) %>%
3   select(variable_name, question_text)
4
5 idea_question_labels %>%
6   gt()

```

variable_name question_text

IDEA_1 Is this period of your life possibilities?

IDEA_2 Is this period of your life

IDEA_3 Is this period of your life out?

IDEA_4 Is this period of your life

IDEA_5 Is this period of your life

IDEA_6 Is this period of your life own beliefs and values?

IDEA_7 Is this period of your life some ways but not others?

IDEA_8 Is this period of your life becoming an adult?

I DEA Questions (VIL)

Because the column contains the left

```
1 idea_responses_clean <- idea_responses_long %>%
2   left_join(idea_questions,
3     by = c("name" = "variable_name"))
4
5 idea_responses_clean %>%
6   head() %>%
7   gt()
```

ResponseID name value question_text

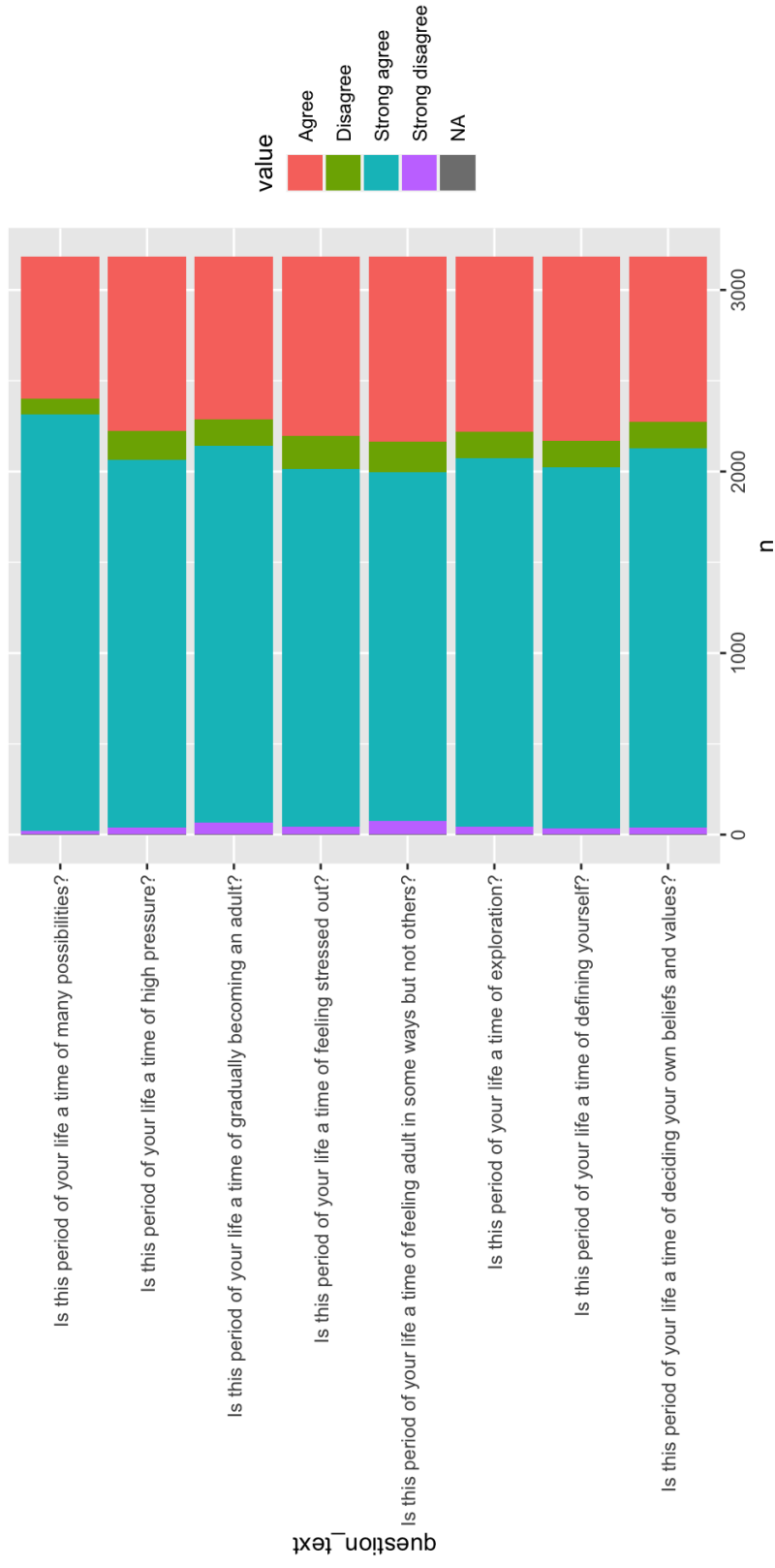
R_BJN3bQqi	1	DEA Agree	Is this period of your life a tir
R_BJN3bQqi	1	DEA Strong	Is this period of your life a tir
R_BJN3bQqi	1	DEA Strong	Is this period of your life a tir
R_BJN3bQqi	1	DEA Agree	Is this period of your life a tir
R_BJN3bQqi	1	DEA Strong	Is this period of your life a tir
R_BJN3bQqi	1	DEA Strong	Is this period of your life a tir

... so why did we do not make a list of the results.

I DE A Q u e s t i o n s (V I I I)

F i n a l l y w e c a n v i s u a l i s e t h e r e s p o n s e s . . .

```
1 idea_responses_clean %>%
2   count(question_text, value) %>%
3   ggplot(aes(x = n,
4             y = question_text,
5             fill = value)) +
6   geom_col()
```

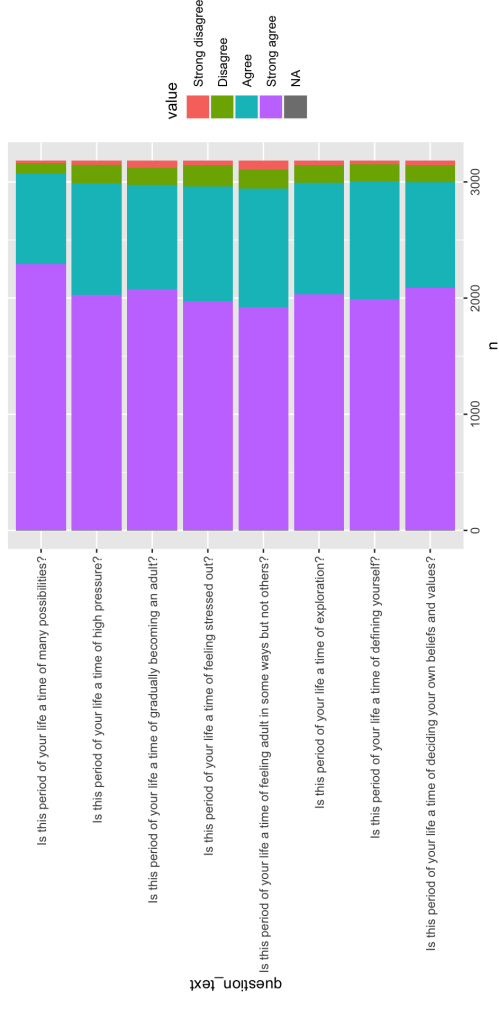


IDEA Questions (IX)

We need to use `fct_relevel()` to set the categories

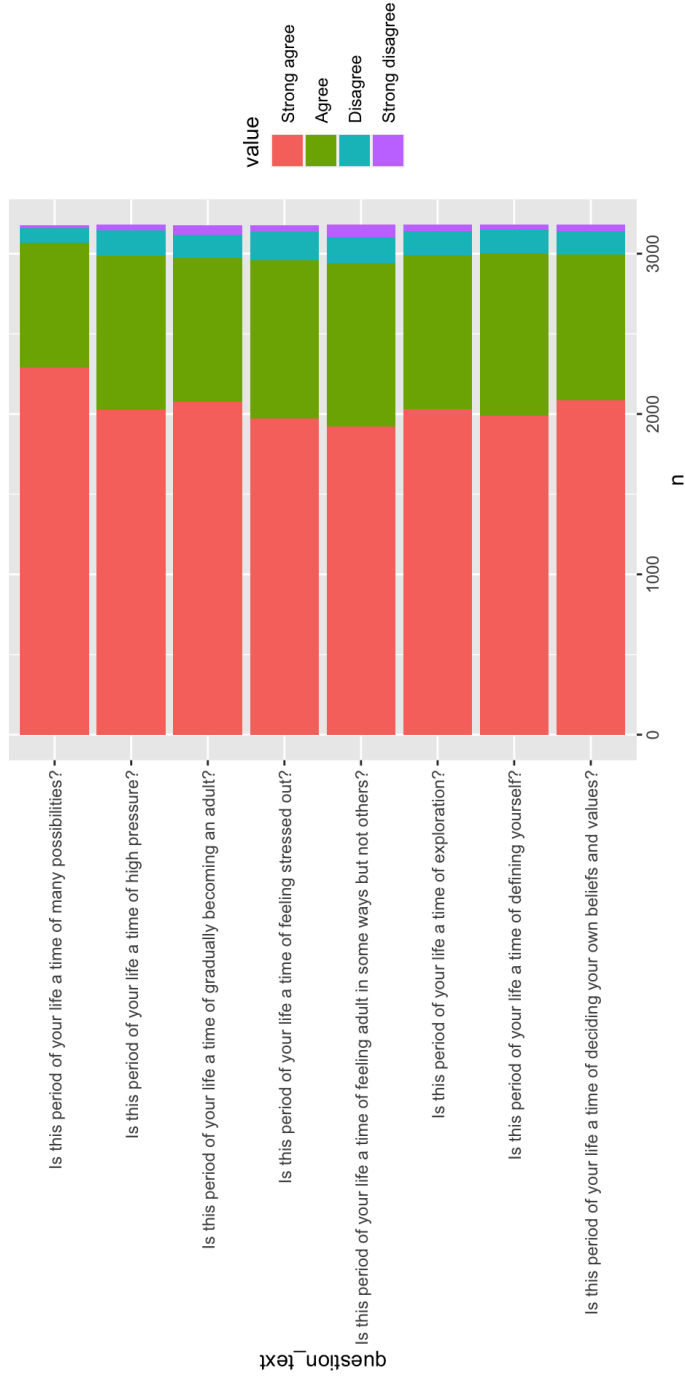
- Which part of the chart do we need to target
- Which part of the chart do we need to target

```
1 order_agree_responses <- c("Strong disagree", "Disagree", "Agree", "Strong agree")
2
3 idea_responses_clean %>%
4   count(question_text, value) %>%
5   mutate(value = fct_relevel(value, order_agree_responses)) %>%
6   ggplot(aes(x = n,
7               y = question_text,
8               fill = value)) +
9   geom_col()
```



I D E A Q u e s t i o n s (X)

```
1 order_agree_responses <- c("Strong disagree", "Disagree", "Agree", "Strong agree")
2
3 idea_responses_clean %>%
4   drop_na(value) %>%
5   count(question_text, value) %>%
6   mutate(value = fct_relevel(value, order_agree_responses)) %>%
7   ggplot(aes(x = n,
8              y = question_text,
9              fill = value)) +
10  geom_col() +
11  scale_fill_discrete(direction = -1) +
12  guides(fill = guide_legend(reverse = TRUE))
```





Exercise: Social m

SLIDE 1 OF 1

Follow this same process to visualise the r
same survey

Lest 'take 20mins for this

Please note that I've not created this char
during tshimsa tvereika'l .



Task: Obtain Commercial Health Data

SLIDE 1 OF 2

1. Register for a FREE UK Data Service account
beta.ukdataservice.ac.uk/myaccount/login
2. Navigate to the access data page for the
beta.ukdataservice.ac.uk/datacatalogue/study
3. Download the SPSS dataset

I st' very useful to learn how to deal with S

1. Unzip the dataset and add the folder to th



Task: Obtain Commercial Health Data

SLIDE 2 OF 2

1. Open up `commercial-access-to-health-data.Rmd`
2. Load `{tidyverse}` package

Read in the dataset (`ipath_data`) really long

```
1 commercial_health_data_raw <- read_spss("data/UKDA-8049-spss_2/spss/spss19/health_data_attitudes_spss_final.sav")
```

Tibbles are great (I)

We've seen before that tibbles are augmented attributes and print more prettily

```
> commercial_health_data_raw
```

```
# A tibble: 2,017 x 124
```

```
mq01_1 mq01_2 mq01_3 mq02_a mq02_b mq02_c mq02_d mq02_e mq03 mq04_1 mq04_2 mq05a mq05b mq06a mq06b mq07_a
<dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl> <dbl+lbl>
1 4 [Heard of, k... 5 [Nev... 5 [Nev... NA NA NA 5 [Agr... 5 [Agr... 5 [Str... 5 [Str... 3 [3.] NA 5 [5. ... 1 [Str... NA 5 [Str... NA mq07_a
2 1 [A great dea... 1 [A g... 1 [A g... NA NA NA 2 [Agr... NA NA NA 1 [Str... NA NA 2 [Ten... 4 [Agr...
3 3 [Just a litt... 2 [A f... 2 [A f... 5 [Agr... NA NA NA NA NA NA 2 [Ten... 2 [Ten... 2 [Ten... 1 [Agr...
4 1 [A great dea... 1 [A g... 1 [A g... NA NA NA NA NA NA 2 [Agr... 3 [Agr... 2 [Ten... 2 [Ten... 2 [Agr...
5 3 [Just a litt... 3 [Jus... 3 [Jus... NA NA NA NA NA NA 2 [Agr... 3 [Agr... 2 [Ten... 2 [Ten... 2 [Agr...
6 3 [Just a litt... 2 [A f... 3 [Jus... NA NA NA NA NA NA 2 [Agr... 3 [Agr... 2 [Ten... 2 [Ten... 2 [Agr...
7 3 [Just a litt... 2 [A f... 2 [A f... NA NA NA NA NA NA 4 [Agr... 3 [Agr... 2 [Ten... 2 [Ten... 4 [Agr...
8 3 [Just a litt... 3 [Jus... 4 [Head... NA NA NA NA NA NA 3 [Agr... 3 [Agr... 2 [Ten... 2 [Ten... 3 [Agr...
9 3 [Just a litt... 3 [Jus... 3 [Jus... NA NA NA NA NA NA 3 [Agr... 2 [Ten... 2 [Ten... 2 [Ten... 2 [Agr...
10 3 [Just a litt... 3 [Jus... 3 [Jus... NA NA NA NA NA NA 6 [Agr... 2 [Ten... 4 [Ten... 5 [5. ... NA 5 [Str... NA
```

```
# ... with 2,007 more rows, and 108 more variables: mq07_b <dbl+lbl>, mq08a1 <dbl+lbl>, mq08a2 <dbl+lbl>, mq08a3 <dbl+lbl>, mq08a4 <dbl+lbl>,
# mq08a5 <dbl+lbl>, mq08a6 <dbl+lbl>, mq08a7 <dbl+lbl>, mq08a8 <dbl+lbl>, mq08a9 <dbl+lbl>, mq08a10 <dbl+lbl>, mq08a11 <dbl+lbl>,
# mq08a12 <dbl+lbl>, mq08b <dbl+lbl>, region <dbl+lbl>, age3 <dbl+lbl>, sex <dbl+lbl>, work <dbl+lbl>, cie <dbl+lbl>, mshop <dbl+lbl>,
# super <dbl+lbl>, wrkcie <dbl+lbl>, sgrade <dbl+lbl>, maritl <dbl+lbl>, numhhd <dbl+lbl>, numkid <dbl+lbl>, numkid2 <dbl+lbl>,
# numkid31 <dbl+lbl>, numkid32 <dbl+lbl>, numkid33 <dbl+lbl>, numkid34 <dbl+lbl>, numkid35 <dbl+lbl>, numkid36 <dbl+lbl>, dura1 <dbl+lbl>,
# dura2 <dbl+lbl>, dura3 <dbl+lbl>, dura4 <dbl+lbl>, dura5 <dbl+lbl>, dura6 <dbl+lbl>, dura7 <dbl+lbl>, dura8 <dbl+lbl>, dura9 <dbl+lbl>,
# dura10 <dbl+lbl>, dura11 <dbl+lbl>, dura12 <dbl+lbl>, dura13 <dbl+lbl>, dura14 <dbl+lbl>, dura15 <dbl+lbl>, dura16 <dbl+lbl>, ...
# i Use `print(n = ...)` to see more rows, and `colnames()` to see all variable names
```

Ti b b l e s a r e g r e a t (l l

The {haven} package currently saves serially confidential "labelled" column class which contains:

- singular "label" means
- The question label (or a shortened version
- The question values • plural "labels" means
- The question value label.

```
> commercial_health_data_raw$mq01_1
<labelled<double>[2017]>: MQ01_1 - Health data collected from patients in hospitals and GP practices can also be used for research
[1] 4 1 3 1 3 3 3 3 4 4 3 1 2 4 4 2 2 4 2 3 4 3 3 2 4 1 1 3 1 1 1 1 1 4 3 1 3 2 4 2 3 2 4 3 2 3 5 5 5 4 4 5 5 6 5 2 5 1 4 3 3 3 4 4 2 4 3 3
[71] 3 3 1 3 3 3 3 3 3 3 3 2 5 2 4 4 4 1 4 4 5 3 5 2 2 4 3 3 5 3 3 2 1 3 5 4 2 2 3 2 4 4 3 5 3 3 3 2 2 3 2 1 1 5 5 3 2 4 3 3
[141] 5 3 4 3 1 4 2 2 3 3 4 4 3 3 1 1 1 1 2 1 4 4 4 5 2 4 4 4 3 4 1 2 1 2 4 1 2 3 4 2 1 1 2 2 4 6 6 3 3 3 2 4 4 1 4 4 3 3 5 3 3 3 5 3
[211] 4 4 3 5 2 4 3 5 3 4 1 2 4 4 2 2 3 3 3 3 1 1 1 1 2 1 4 4 4 5 2 4 4 4 3 4 1 2 1 2 4 1 2 3 4 2 1 1 2 2 4 6 6 3 3 3 2 4 4 1 4 4 3 3 5 3 3
[281] 3 5 1 2 3 1 5 4 4 4 3 2 4 3 1 3 2 4 2 1 3 3 5 4 2 5 4 4 5 5 3 5 2 4 3 5 3 3 4 5 4 1 2 2 2 3 2 2 3 3 1 2 4 2 1 2 1 3 3 3 2 2 2 3 2
[351] 5 5 3 3 5 5 5 4 4 3 5 4 5 4 4 1 3 3 2 4 4 3 3 2 2 5 3 3 5 3 3 2 4 3 3 3 1 3 3 4 1 2 1 3 5 5 2 1 5 2 2 3 2 5 2 5 2 2 2 1 2 1
[421] 4 5 3 4 3 4 3 2 3 1 4 4 2 2 2 2 3 3 5 3 5 2 1 2 4 3 2 3 5 3 3 5 3 3 5 3 3 2 3 5 4 3 1 5 5 3 2 3 2 4 2 3 4 5 3 5 4 3 5 3 2 3 3 4 5 2 2
[491] 2 3 4 3 3 2 3 1 4 4 2 5 3 2 2 2 3 2 2 4 2 3 5 3 3 2 3 5 2 2 3 1 4 3 3 2 3 3 1 3 3 6 2 5 3 4 4 3 2 5 3 1 2 5 4 3 5 4 2 5 1 4 4
[561] 3 2 4 5 2 3 3 2 6 3 3 2 3 2 5 3 3 3 4 1 1 1 3 2 1 6 2 2 2 1 5 5 4 1 2 3 2 2 2 2 2 3 4 3 2 2 3 4 3 4 5 5 5 1 2 2 5 3 4 4 4 1
[631] 3 3 3 5 5 5 2 5 3 3 5 3 4 4 4 4 3 4 3 5 4 4 4 4 3 4 4 5 3 1 2 1 1 2 2 3 5 5 4 4 5 4 4 5 4 2 5 2 6 5 4 1 4 4 2 1 2 4 4 2 1 3 2 2 3
[701] 1 4 1 3 1 3 1 4 3 5 2 4 1 2 4 4 3 4 5 3 5 4 1 3 1 1 1 4 5 3 3 5 1 1 1 2 4 5 2 4 4 4 3 2 3 4 4 3 3 3 2 4 2 3 1 2 3 2 1 3 5 2 3 4
[771] 3 2 2 2 5 2 4 3 3 2 2 3 2 2 5 6 3 3 4 4 2 3 1 3 2 2 4 4 1 5 5 5 2 3 3 3 3 3 3 5 3 5 1 5 5 1 1 3 2 3 5 3 4 4 1 1 1 1 2 1 4
[841] 3 1 3 2 3 4 4 3 4 3 4 4 1 3 1 4 3 2 1 3 4 1 3 4 2 3 1 2 2 1 4 3 2 2 3 5 3 1 4 3 3 4 5 4 3 4 2 5 4 4 5 4 5 3 5 3 2 2 5 2 5 2 3
[911] 5 5 1 4 4 5 2 4 3 4 4 2 5 4 4 2 5 3 2 5 1 4 4 2 1 4 4 2 2 2 5 3 2 4 2 2 3 2 2 4 3 3 5 2 4 2 4 5 2 4 5 4 5 2 4 5 2 4 3 5 3 4 3
[981] 3 2 3 3 3 1 4 5 3 3 4 1 4 5 3 3 2 3 3
[ reached getOption("max.print") -- omitted 1017 entries ]
```

```
Labels:
value      label
1          A great deal
2          A fair amount
3          Just a little
4  Heard of, know nothing about
5  Never heard of
6  Don't know
7  MISSING
```

Extracting the quest {haven} output (I)

We have two ways to extract labels from {

- Programming with purrr::attr_getter()

```
1 commercial_health_data_raw %>%  
2   map_chr(attr_getter("label"))
```

- Using sjlabelled::get_label()

```
1 sjlabelled::get_label(commercial_health_data_raw)
```

Unfortunately programming solution only works

I'd recommend [sjlabelled](#) package exclusively for labels from SPSS. It does provide [functions](#) for also calculating

Extracting the question {haven} output (II)

The `sjlabelled`: `get_f_u_n_c_t_i_l` `q_n` generated answers
into a `teinbraemew(i)` `th`

```
1 commercial_health_data_qs_raw <- sjlabelled(:get_label(commercial_health_data_raw) %>%
2   enframe() %>%
3   rename(question_text = value)
4 commercial_health_data_qs_raw
```

```
# A tibble: 124 x 2
  name      question_text
<chr>   <chr>
1 mq01_1 MD01_1 - Health data collected from patients in hospitals and GP prac...
2 mq01_2 MD01_2 - Health data collected from patients in hospitals and GP prac...
3 mq01_3 MD01_3 - Health data collected from patients in hospitals and GP prac...
4 mq02_a MD02_A - As you may know the NHS and other health services collect d...
5 mq02_b MD02_B - As you may know the NHS and other health services collect d...
6 mq02_c MD02_C - As you may know the NHS and other health services collect d...
7 mq02_d MD02_D - As you may know the NHS and other health services collect d...
8 mq02_e MD02_E - As you may know the NHS and other health services collect d...
9 mq03   MD03 - To what extent, if at all, would you support your health data ...
10 mq04_1 MD04_1 - To what extent do you agree or disagree with the following s...
# ... with 114 more rows
```

Can you help me write some code to remove
column?

Extracting the question {haven} output (II)

This is one of many ways to tidy up this data

```
1 commercial_health_data_qs <- commercial_health_data_qs_raw %>%
2   mutate(question_text = str_remove(question_text, toupper(name)),
3          question_text = str_remove(question_text, " - "),
4          question_text = str_remove(question_text, "MO08A"))
5 commercial_health_data_qs
```

```
# A tibble: 124 × 2
   name  question_text
<chr> <chr>
1 mq01_1 Health data collected from patients in hospitals and GP practices can...
2 mq01_2 Health data collected from patients in hospitals and GP practices can...
3 mq01_3 Health data collected from patients in hospitals and GP practices can...
4 mq02_a As you may know the NHS and other health services collect data about...
5 mq02_b As you may know the NHS and other health services collect data about...
6 mq02_c As you may know the NHS and other health services collect data about...
7 mq02_d As you may know the NHS and other health services collect data about...
8 mq02_e As you may know the NHS and other health services collect data about...
9 mq03   To what extent, if at all, would you support your health data being a...
10 mq04_1 To what extent do you agree or disagree with the following statements?
# ...with 114 more rows
```

Converting labelled

To convert all labelled columns to factors

```

1 commercial_heal_th_data_factors <- commercial_heal_th_data_raw %>%
2   mutate(across(where(is.labelled), ~as_factor(.x)))
3 commercial_heal_th_data_factors

# A tibble: 2,017 × 124
  mq01_1 mq01_2 mq01_3 mq02_a mq02_b mq02_c mq02_d mq02_e mq03 mq04_1 mq04_2
  <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct>
1 Heard o... Never ... Never ... <NA> <NA> <NA> Agree... Stro... Stro... Stro... Stro...
2 A great ... A gre... A gre... <NA> Agree... <NA> <NA> <NA> Stro... Stro... Stro... Stro...
3 Just a ... A fai ... A fai ... Agree... <NA> <NA> <NA> Tend... Tend... Tend... Tend...
4 A great ... A gre... A gre... <NA> <NA> <NA> <NA> Agree... Nei t... Tend ... Tend ...
5 Just a ... Just ... Just ... <NA> <NA> <NA> <NA> Agree... Tend... Nei th... Tend ...
6 Just a ... A fai ... Just ... <NA> <NA> <NA> <NA> Agree... Tend... Tend ... Tend ...
7 Just a ... A fai ... A fai ... <NA> <NA> <NA> <NA> Agree... Nei t... Tend ... Tend ...
8 Just a ... Just ... Heard... <NA> <NA> <NA> <NA> Agree... Nei t... Nei th... Tend ...
9 Just a ... Just ... Just ... <NA> <NA> <NA> <NA> Agree... Tend... Tend ... Tend ...
10 Just a ... Just ... Just ... <NA> <NA> Agree... <NA> Tend... Tend ... Stro...
# ... with 2,007 more rows, and 113 more variables: mq05a <fct>, mq05b <fct>,
# mq06a <fct>, mq06b <fct>, mq07_a <fct>, mq07_b <fct>, mq08a1 <fct>,
# mq08a2 <fct>, mq08a3 <fct>, mq08a4 <fct>, mq08a5 <fct>, mq08a6 <fct>,
# mq08a7 <fct>, mq08a8 <fct>, mq08a9 <fct>, mq08a10 <fct>, mq08a11 <fct>,
# mq08a12 <fct>, mq08b <fct>, region <fct>, age3 <fct>, sex <fct>,
# work <fct>, cie <fct>, mshop <fct>, super <fct>, wrkcie <fct>,
# sgrade <fct>, maritl <fct>, numhhd <fct>, numkid <fct>, numki d2 <fct>, ...

```

Notice how we don't have a respondent ID

Add respondent ID

The `row_number(n)` function gives us a neat way to add a new column to a data frame. However, it's not necessarily that clever a solution.

```
1 commercial_health_data_clean <- commercial_health_data_factors %>%
2   mutate(respondent_id = row_number()) %>%
3   relocate(respondent_id)
4 commercial_health_data_clean
```

```
# A tibble: 2,017 x 125
  respondent_id mq01_1 mq01_2 mq01_3 mq02_a mq02_b mq02_c mq02_d mq02_e mq03 mq04_1
    <int> <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct>
1      1 Heard... Never ... <NA> <NA> <NA> Agree... Stro... Stro...
2      2 A gre... A gre... <NA> Agree... <NA> <NA> Stro... Nei th...
3      3 Just ... A fai ... A fai ... Agree... <NA> <NA> Tend... Tend ...
4      4 A gre... A gre... <NA> <NA> <NA> Agree... Nei t... Tend ...
5      5 Just ... Just ... Just ... <NA> <NA> Agree... Tend... Nei th...
6      6 Just ... A fai ... Just ... <NA> <NA> Agree... Tend... Tend ...
7      7 Just ... A fai ... A fai ... <NA> <NA> Agree... Nei t... Tend ...
8      8 Just ... Just ... Heard... <NA> <NA> Agree... Nei t... Nei th...
9      9 Just ... Just ... Just ... <NA> <NA> Agree... Tend... Tend ...
10     10 Just ... Just ... Just ... <NA> Agree... <NA> <NA> Tend... Tend ...
# ... with 2,007 more rows, 114 more variables: mq04_2 <fct>, mq05a <fct>,
# mq05b <fct>, mq06a <fct>, mq06b <fct>, mq07_a <fct>, mq07_b <fct>,
# mq08a1 <fct>, mq08a2 <fct>, mq08a3 <fct>, mq08a4 <fct>, mq08a5 <fct>,
# mq08a6 <fct>, mq08a7 <fct>, mq08a8 <fct>, mq08a9 <fct>, mq08a10 <fct>,
# mq08a11 <fct>, mq08a12 <fct>, mq08b <fct>, region <fct>, age3 <fct>,
# sex <fct>, work <fct>, cie <fct>, nshop <fct>, super <fct>, wrkie <fct>,
# sgrade <fct>, mariti <fct>, numrhd <fct>, numkid <fct>, numkid2 <fct>, ...
```

Commercial Health Data

I'd like you to extract the columns from the

Q4. To what extent do you agree or disagree with the following statements?

	"My health data currently has financial value to others in that it can be used to save or make them money."	"My health data currently has a value to society in that it can be used to help improve things for people other than me."
Base:	<i>All respondents (2,017)</i> %	<i>All respondents (2,017)</i> %
Strongly agree	15	28
Tend to agree	35	40
Neither agree nor disagree	25	18
Tend to disagree	12	7
Strongly disagree	9	5
Don't know	3	3
Agree	50	67
Disagree	21	12

Commercial Health Data

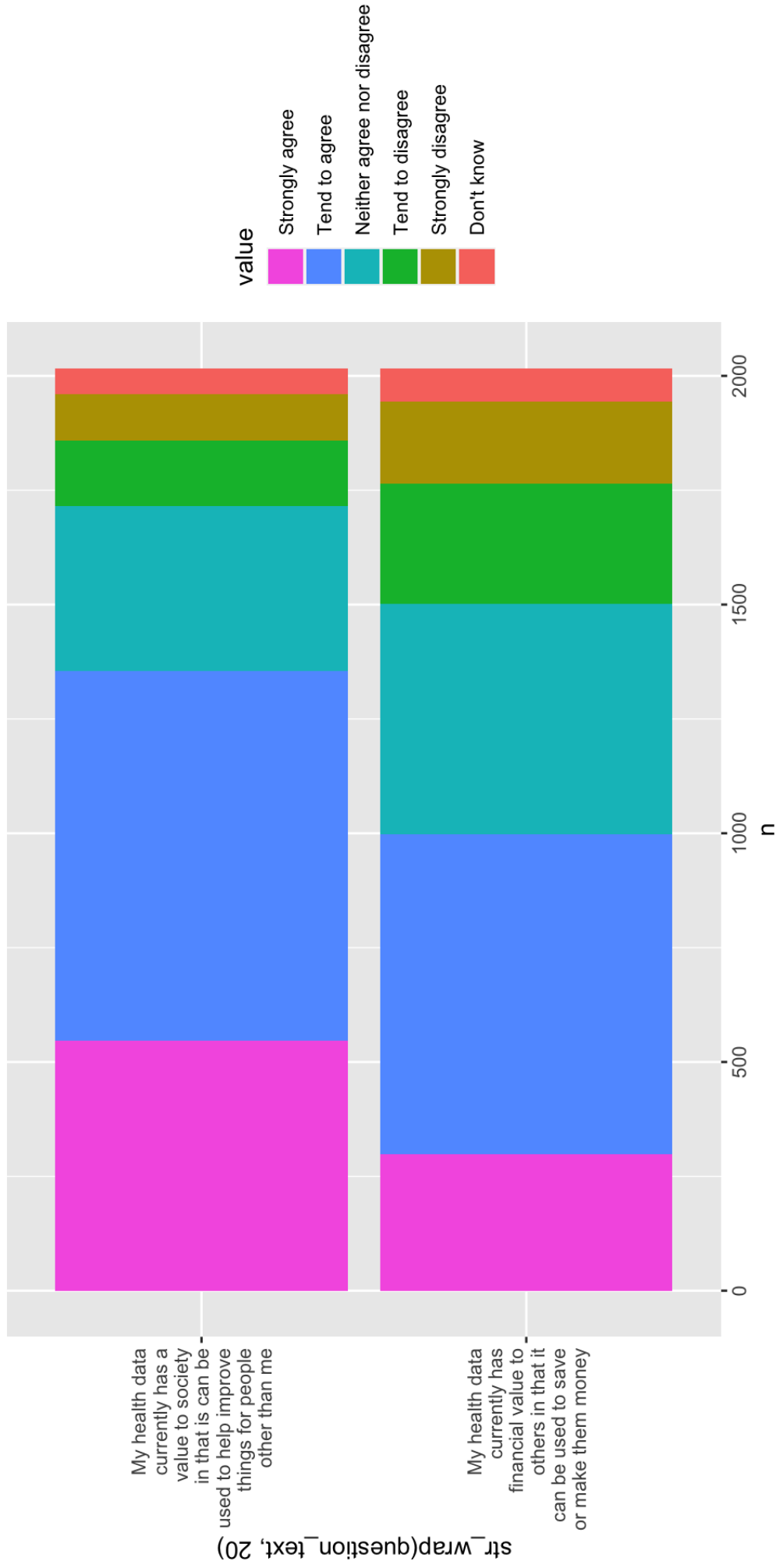
What do we need to do this data so that
{ ggplot 2 }

```
1 commercial_health_data_clean %>%  
2 select(respondent_id, starts_with("mq04"))
```

```
# Attibble: 2,017 x 3  
  respondent_id mq04_1  
    <int> <fct>  
1      1 Strongly disagree  
2      2 Neither agree nor disagree  
3      3 Tend to agree  
4      4 Tend to agree  
5      5 Neither agree nor disagree  
6      6 Tend to agree  
7      7 Tend to agree  
8      8 Neither agree nor disagree  
9      9 Tend to agree  
10     10 Tend to disagree  
# ... with 2,007 more rows  
mq04_2  
    <fct>  
1 Strongly disagree  
2 Strongly agree  
3 Tend to agree  
4 Tend to agree  
5 Tend to agree  
6 Tend to agree  
7 Tend to agree  
8 Tend to agree  
9 Tend to agree  
10 Strongly disagree
```

Commercial Health Data

Let's make this chart:



How wide is too


```
1  tweetrmd: : tweet_embed("https://twitter.com/charliejhadley/status/1522559488284413954?ref_src=twsrc%5Etfw")
```

Buffy ratings

In what ways is this dataset wide?

- Ratings are split across multiple columns
- Should we [join](#) [dream](#) [ratings](#)?
- In principle we could combine:

- votes
- views
- `voxep_rank`

```
1 buffy_raw <- read_csv(here::here("static", "datasets", "buffy", "buffy_data.csv"))
2 buffy_raw
```

```
# A tibble: 144 × 15
  no overall ..1 season no in..2 title direc...3 writer air_d... views... imdb_... votes
  <dbl> <dbl> <dbl> <dbl> <chr> <chr> <chr> <chr> <dbl> <dbl> <dbl> <dbl>
1 1 1 1 1 1 W...c... Charles... Joss ... 3/10/1... <chr> <dbl> 8 4548
2 2 1 1 1 2 The ... John T... Joss ... 3/10/1... <chr> <dbl> 7.8 3952
3 3 1 1 1 3 W...tch Stephe... Dana ... 3/17/1... <chr> <dbl> 7.7 3940
4 4 1 1 1 4 Teac... Bruce ... David... 3/24/1... <chr> <dbl> 6.9 3800
5 5 1 1 1 5 Neve... David ... Rob D... 3/31/1... <chr> <dbl> 7.4 3611
6 6 1 1 1 6 The ... Bruce ... Matt ... 4/7/19... <chr> <dbl> 7.4 3770
7 7 1 1 1 7 Angel Scott ... David... 4/14/1... <chr> <dbl> 8.5 3918
8 8 1 1 1 8 I, R... Stephe... Ashle... 4/28/1... <chr> <dbl> 6.7 3629
9 9 1 1 1 9 The ... Ellen ... Rob D... 5/5/19... <chr> <dbl> 7.7 3666
10 10 1 1 1 10 Nigh... Bruce ... Joss ... 5/12/1... <chr> <dbl> 8.2 3614
# ... with 134 more rows, 5 more variables: plot <chr>, runtime <dbl>,
# deathcount <dbl>, `neilsen rating` <chr>, `voxep_rank` <dbl>, and
```

Long enough for what

Tidy data is a useful concept for ⁴wrangling

Buts into't something to. conform to religiousl

You might want to keep some width to your d

Wide data might also be more appropriate if

Other forms of un

Multiplying pieces of data

Sometimes a single column might ask for multiple variables. all that apply"

This is often the case when you live?"

```
1 location_data <- tribble(
2   ~id, ~address,
3   1, "Las Vegas, USA",
4   2, "Bristol, UK",
5   3, "Kassala, Sudan"
6 )
7 location_data
```

```
# A tibble: 3 × 2
  id address
<dbl> <chr>
1 1 Las Vegas, USA
2 2 Bristol, UK
3 3 Kassala, Sudan
```

```
1 device_ownership <- tribble(
2   ~name, ~devices_owned,
3   "Charlie", "Smart TV, Cell phone",
4   "Mohammad", "Cell phone",
5   "Christina", "Smart TV, Games Console, Cell phone"
6 )
7 device_ownership
```

```
# A tibble: 3 × 2
  name devices_owned
<chr> <chr>
1 Charlie Smart TV, Cell phone
2 Mohammad Cell phone
3 Christina Smart TV, Games Console, Cell phone
```



Task: Obtain British Data

SLIDE 1 OF 2

1. Register for a FREE British [Ecological Society
login.php?action=register](https://www.britishecologicalsociety.org/researchers/register)
2. Navigate to the access [British
object/2019-british-election-study-post-el](https://www.britishecologicalsociety.org/researchers/register)
3. Download the SPSS dataset
4. Unzip the dataset and add the folder to the



Task: Obtain British Data

SLIDE 2 OF 2

1. Setup the British. Remove identifiers as well
2. Read in the SPSS

```
# A tibble: 3,946 x 415
  final ser... agency Y10A Y10B1 Y10B2 Y10B3 Y10B4 Y10B5 a01
  <dbl> <dbl> +l b> <dbl> <dbl> +l b> <dbl> +l b> <dbl> +l b> <dbl> +l b> <chr> +l bl >
1 10102 1 [l ps... NA NA 0 [No] 0 [No] 0 [No] 0 [No] ""
2 10103 NA 2 1 [Yes] 0 [No] 0 [No] 0 [No] 0 [No] ""
3 10105 NA 2 1 [Yes] 0 [No] 0 [No] 0 [No] 0 [No] ""
4 10110 1 [l ps... NA NA NA NA NA "- 1" [Don...
5 10111 1 [l ps... NA NA NA NA NA "- 1" [Don...
6 10202 NA 2 1 [Yes] 0 [No] 0 [No] 0 [No] 0 [No] ""
7 10206 NA 3 1 [Yes] 0 [No] 0 [No] 0 [No] 0 [No] ""
8 10208 NA 2 1 [Yes] 0 [No] 0 [No] 0 [No] 0 [No] ""
9 10210 NA 2 1 [Yes] 0 [No] 0 [No] 0 [No] 0 [No] ""
10 10304 NA 2 1 [Yes] 0 [No] 0 [No] 0 [No] 0 [No] ""
# ... with 3,936 more rows, 406 more variables: a01_code <dbl> +l bl >,
# a02 <dbl> +l bl >, a03 <dbl> +l bl >, m02_1 <dbl> +l bl >, m02_2 <dbl> +l bl >,
# m02_3 <dbl> +l bl >, m02_4 <dbl> +l bl >, m02_5 <dbl> +l bl >, m02_6 <dbl> +l bl >,
# b01 <dbl> +l bl >, b02 <dbl> +l bl >, b04 <dbl> +l bl >, b05 <dbl> +l bl >,
# b0601 <dbl> +l bl >, b0602 <dbl> +l bl >, b0603 <dbl> +l bl >, b0604 <dbl> +l bl >,
# b0605 <dbl> +l bl >, b0606 <dbl> +l bl >, b0607 <dbl> +l bl >, b0608 <dbl> +l bl >,
# b0609 <dbl> +l bl >, b0610 <dbl> +l bl >, b0611 <dbl> +l bl >, b0612 <dbl> +l bl >, ...
```

Where do people get from?

Can you extract the corresponding

dataset corresponding to each country?

KO4: Where do you get most of your information about politics or current affairs from? (Modes: CAP/Online/Paper. Countries: England/Scotland/Wales.)

Value	Label
-1	Don't know
-2	Refused

What can you tell me about this data
and question?

Where do people get from?

This is an open-ended
going to be really messy

```
1 british_election_data_raw %>%  
2 select(final_serial_no, k04)
```

A tibble: 3,946 × 2
final_serial_no k04

```
1 10102 "family"  
2 10103 "-2" [Refused]  
3 10105 "Media- cross referencing and watching  
parliamentary debates"  
4 10110 "parents"  
5 10111 "tv radio"  
6 10202 ""  
7 10206 "News, internet and conversation"  
8 10208 ""  
9 10210 "Mail on line \nNews on tv"  
10 10304 "tv papers."  
# ... with 3,936 more rows
```

To properly analyse the tidytext package
with a tidyverse approach

But still, what we can do by
pretending it's multiple choice data and
using
separate()
separate_rows()

References

1. Gri n N. et al The i portance of urban natural areas and urban ecosystem servi es duri g the COVID-19 pandemi .c PLOS ONE 15, e0243344 (2020).
2. Eatock, J., Di on, D. & Young, T. An expl ratory survey of current practi ei the medi at devi ei dustry. Journal of Manufacturi g Technol gy Management 20, 218– 234 (2009).
3. Presser, S. et al Methods for Testi grand Eval ai g Survey Questi es. Publ iOpi inoQuarterl y8, 109–130 (2004).
4. Wi kham, H. Ti yData. Journal of Stati ts at Software 59, 1– 23 (2014).
5. Fi lehouse, E. et al Bri i tsEI ei roStudy, 2019: Post-EI ei roRandom Probabi ty Survey. (2019) doi 10.52555/UKDA-SN-8875-1.
6. NHS Engl and. Wri i tti an effecti e questi onnai e. (2018).
7. Gal p. Why Phone and Web Survey Resul stAren' the Same. Gal p.com (2018).
8. Sánchez Tomé, R. The i pact of mode of data col lti on measures of subj ei ewel ei bti (Uni ersi ty of Lausanne, 2018).
9. AAPOR. An Eval ai rof 2016 EI ei roPol li sthe U.S. Ameri an Associ ti rofor Publ iOpi inoResearch (2016).
10. Serdar, C. C., Ci an, M., Yücel D. & Serdar, M. A. Sampl ei e power and effecti e e revi ist: Si plified and practi at approaches i pre-cl i al cl i ad and l boratory studi e Bi che mi Med i al31, 010502 (2021).
11. Perneger, T. V., Courvoi isre D. S., Hudel an, P. M. & Gayet-Ageron, A. Sampl ei e for pre-tests of questi onnai es. Qual y of Li eResearch 24, 147–151 (2015).
12. Tel l, S. Sampl ei e How many peopl e shoul dake the survey? on devi e research (2014).
13. Jami an, S. Li rt scal e How to (ab)use them. Medi al Educati o38, 1217–1218 (2004).
14. Batterton, K. A. & Hal , K. N. The Li rt Scal e What It Is and How To Use It. Phal x50, 32– 39 (2017).
15. SurveyMonkey. Does addi g more questi oi pact survey compl te rate? SurveyMonkey.
16. Rhemtul l M., Saval ieV. & Li tt, et. D. On the Asymptoti e el ti e/Effici ecy of PI aned Mi s gness Desi es.
16. Psychometri al81, 60–89 (2016).
17. Grahe, J. E. et al Emergi gAdul hoo Measured at Mul i tti onsti uti es 2: The Data. Journal of Open Psychol gy Data 6, 4 (2018).
18. Reifan, A. & Grahe, J. E. Introducti ro to the Speci alssue of Emergi gAdul hoo. Emergi gAdul hoo 4, 135–141 (2016).
19. Grahe, J. et al EAMMi e Publ iData. (2022) doi 10.17605/OSF.IO/QTQPB.