

Week 3: Data science tidyverse

Charlotte Hadley

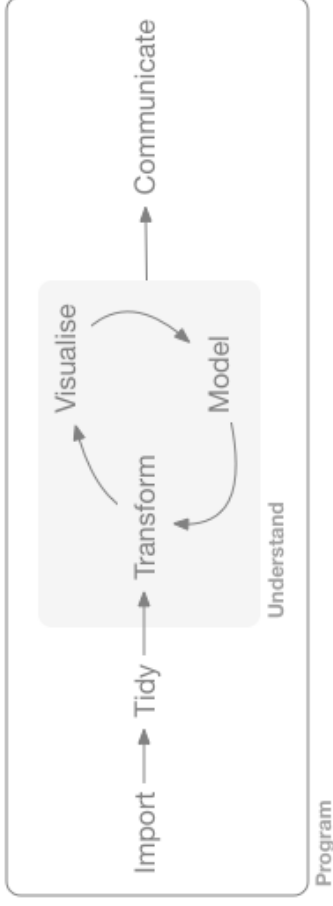
Topics for today

1. Understanding a `dataframe`
2. Using `pandas` for reading data
3. Using `pandas` for cleaning and wrangling data
4. Using `ggplot2` for some basic exploratory data

What is the tidy

The tidyverse is two

- A collection of packages designed to work together
- A collection of packages for data science workflow



Source: [Data Science](#)

After we've been using the tidyverse for a while, we can provide further context.

- An opinionated framework for working with data.

Some people would describe the tidyverse as being an R base R.

It's not.

The tidyverse is a useful working with data that ecosystem of libraries. It works quickly.



Task: Create a week-

SLIDE 1 OF 1

1. Create a new RStudio project called something

Installing and working with tidyverse

To install the tidyverse collection of packages

```
1 install.packages("tidyverse")
```

The tidyverse packages are split into two groups

- Core tidyverse packages published "package
loaded with this code explicitly read from
importing data from E

```
1 library(tidyverse)
```

```
1 library(tidyverse)  
2 library(readxl)
```

Updating the tidyverse

In terms of real-world usage, keeping the tidyverse packages up to date.

Thee

- When installing a package the console might prompt you to update all of the packages

These packages have more recent versions available.
It is recommended to update all of them.
Which would you like to update?

```
1: All
2: CRAN packages only
3: None
4: viridisLite (0.4.0 -> 0.4.1) [CRAN]
```

Enter one or more numbers, or an empty line to skip updates: |

Updating the tidyverse

In terms of real-world usage, keeping the tidyverse package up to date.

There are three ways you might discover you should update:

- When installing a package the console might prompt you to update the package.
- When installing a package the install fails due to a dependency issue.
- You hear about an exciting new update to a package.

There is a way to keep the tidyverse up to date automatically, but it is not very rarely used.

Data sets we'll be using

We're going to be using data sets today:

- The Global Burden of Disease Study³. Data from Ethn

The Global Burden of Disease study is an excellent understanding global (and comparative) health

These excellent interactive tool for building a need to register for a free account to use

Data sets we'll be using

We're going to be using datasets today.

- The Global Burden of Diseases³.
- The `msleep` dataset from `gap` package

Lots (and lots) of R packages built into the demonstration how to use the package.

The `msleep` dataset has data on mammalian sleep cycles. We⁴st

1 `glimpse(msleep)`

Rows: 83

Columns: 11

```
$ name      <chr> "Cheetah", "Owl monkey", "Mountain  
$ genus     <chr> "Acinonyx", "Aotus", "Aplodontia",  
"Blarina", "Bos", "Bracon",  
$ vore      <chr> "carni", "omi", "herbi", "omi",  
"herbi", "herbi", "carni",  
$ order     <chr> "Carnivora", "Primates",  
"Rodentia", "Soricomorpha", "Art...  
$ conservation <chr> "lc", NA, "nt", "lc",  
"domesticated", NA, "vu", NA, "dome...  
$ sleep_total <dbl> 12.1, 17.0, 14.4, 14.9, 4.0, 14.4,  
8.7, 7.0, 10.1, 3.0, 5...  
$ sleep_rem  <dbl> NA, 1.8, 2.4, 2.3, 0.7, 2.2, 1.4,  
NA, 2.9, NA, 0.6, 0.8, ...  
$ sleep_cycle <dbl> NA, NA, 0.1333333, 0.6666667,  
0.7666667, 0.3833333, N...  
$ awake     <dbl> 11.9, 7.0, 9.6, 9.1, 20.0, 9.6,  
15.3, 17.0, 13.9, 21.0, 1...  
$ brainwt   <dbl> NA, 0.01550, NA, 0.00029, 0.42300,  
NA, NA, NA, 0.07000, O...  
$ bodywt    <dbl> 50.000, 0.480, 1.350, 0.019,  
400.000, 3.850, 20.400, 0.04
```

Data sets we'll be using

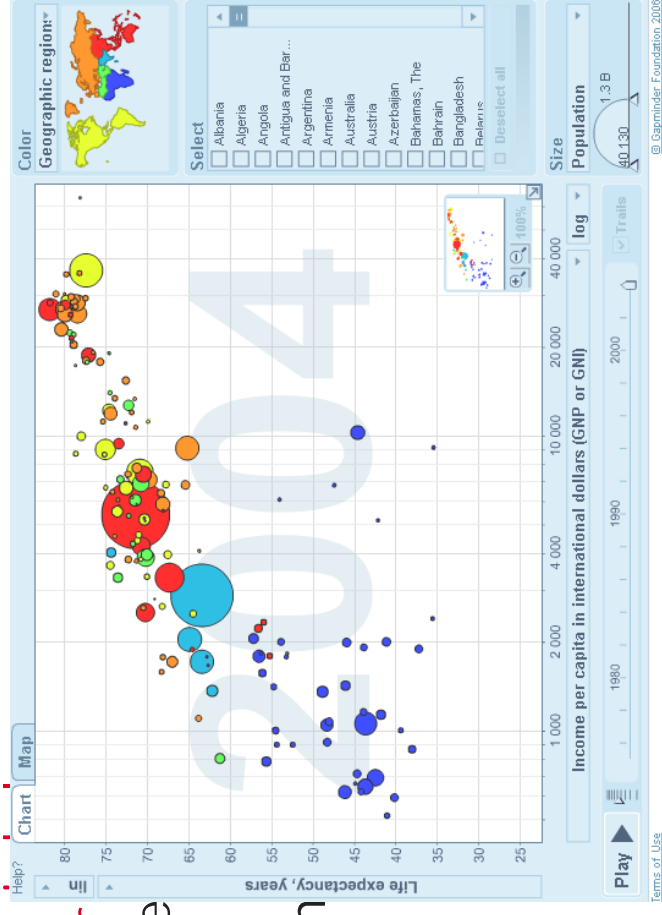
We're going to be using data that also tells the story:

- The Global Burden of **Disability**.
- The **World Development Indicators**.
- The **Gapminder** dataset of **world development**.

In 2006 Hans Rosling **emerged** on the TED stage where he introduced a **simple** **visual** tool to the global development.

Hans Rosling also founded the **Gapminder** Foundation to promote development.

The **Gapminder** dataset of their data.



msl eep



Task: Get the mslleep

SLIDE 1 OF 1

1. Add a heading for the mslleep dataset.
2. Load `tidyverse` package in the setup code chunk
3. Add a new code chunk `mslleep` to import dataset

1 mslleep

Understanding our ob

The datasets `embeddings` and `those` genres
datafiles with `{reader}`, `{readxl}` and `{haven}` a
It's important to understand how to transform

data.frames (I)

The `data.frame` general purpose structure `{qqgjprleadtch}` or build

The `data.sepackage` has a `dambefr`, a `ffedbr` in its ann c

```
1 head(quakes)
```

	lat	long	depth	mag	stations
1 -	20.42	181.62	562	4.8	41
2 -	20.62	181.03	650	4.2	15
3 -	26.00	184.10	42	5.4	43
4 -	17.97	181.66	626	4.1	19
5 -	20.42	181.96	649	4.0	11
6 -	19.68	184.31	195	4.0	12

The `lass(f)` function is the way to `tho/ontogje` `edyædru` in `nee` `wæhræ` with:

```
1 class(quakes)
```

```
[1] "data.frame"
```

data.frames (||)

Because they', defect arrows and columns separately

We can print the data.frames() of a

```
1 dim(quakes)
```

```
[1] 1000 5
```

In Base R there are two ways to extract col

- The \$ operator allows us to extract all columns via their name with both rows and autocompletion: by their index

```
1 quakes$mag
```

```
1 quakes[, "mag"]
```

```
1 quakes[, 1]
```


data.frame is dead, long

The tidyverse ~~data.frame~~ is dead, long live the tidyverse

Let's demonstrate the difference between the two

```
library("tidyverse")
```

The first thing we notice is that the data is not a data frame

```
1 quakes
```

```
1      lat long depth mag stations
2 -20.42 181.62 562 4.8 41
3 -20.62 181.03 650 4.2 15
4 -26.00 184.10 42 5.4 43
5 -17.97 181.66 626 4.1 19
6 -20.42 181.96 649 4.0 11
7 -19.68 184.31 195 4.0 12
8 -11.70 166.10 82 4.8 43
9 -28.11 181.93 194 4.4 15
10 -28.74 181.74 211 4.7 35
11 -17.47 179.59 622 4.3 19
12 -21.44 180.69 583 4.4 13
13 -12.26 167.00 249 4.6 16
14 -18.54 182.11 554 4.4 19
15 -21.00 181.66 600 4.4 10
16 -20.70 169.92 139 6.1 94
17 -15.94 184.95 306 4.3 11
18 -13.64 165.96 50 6.0 83
19 -17.83 181.50 590 4.5 21
20 -23.50 179.78 570 4.4 13
21 -22.63 180.31 598 4.4 18
22 -20.84 181.16 576 4.5 17
23 -10.98 166.32 211 4.2 12
24 -23.30 180.16 512 4.4 18
```

```
1 starwars
```

```
# A tibble: 87 x 14
  name      height mass hair_color skin_color eye_color
  birth... sex      gender homeworld
  <chr>      <chr> <chr> <dbl> <chr> <chr> <chr>
1 Luke Skywalker 172 77 blond fair blue
2 C-3PO          167 75 <NA> gold yellow
3 R2-D2          96 32 <NA> white, ... red
4 Darth Vader    202 136 none white yellow
5 Leia Organa    150 49 brown light brown
6 Owen Lars      178 120 brown, ... light blue
7 Beru Whitesun 165 75 brown light blue
8 R5-D4          97 32 <NA> white, ... red
9 Biggs Darklighter 183 84 black light brown
10 Chirrup         182 77 auburn fair blue-gray
```

data.frame is dead, long

In this course we've explored `readr` and `tidyverse` for data manipulation

The `display` provides a preview of the data

```
1 select(quakes, mag)
```

```
1      4.8
2      4.2
3      5.4
4      4.1
5      4.0
6      4.0
7      4.8
8      4.4
9      4.7
10     4.3
11     4.4
12     4.6
13     4.4
14     4.4
15     6.1
16     4.3
17     6.0
18     4.5
19     4.4
20     4.4
21     4.5
22     4.2
23     4.4
```

```
1 select(starwars, name)
```

```
# A tibble: 87 x 1
  name
<chr>
1 Luke Skywalker
2 C-3PO
3 R2-D2
4 Darth Vader
5 Leia Organa
6 Owen Lars
7 Beru Whitesun Lars
8 R5-D4
9 Biggs Darklighter
10 Obi-Wan Kenobi
# ... with 77 more rows
```

data.frame is dead, long

The select function is not a vector, but some needs a vector

A vector is a one-dimensional atomic object
function:

```
1 c(1, "2", 3)
```

[1] "1" "2" "3"

If we want to extract a value from a vector, use

```
1 pull(starvars, name)
```

- | | | |
|------------------------------|-------------------------|---------------------|
| [1] "Luke Skywalker" | "C-3PO" | "R2-D2" |
| [4] "Darth Vader" | "Leia Organa" | "Owen Lars" |
| [7] "Beru Whitesun Lars" | "R5-D4" | "Biggs Darklighter" |
| [10] "Obi-Wan Kenobi" | "Anakin Skywalker" | "Wilhuff Tarkin" |
| [13] "Chewbacca" | "Han Solo" | "Greedo" |
| [16] "Jabba Desilijic Tiure" | "Wedge Antilles" | "Jek Tono Porkins" |
| [19] "Yoda" | "Palpatine" | "Boba Fett" |
| [22] "IG-88" | "Bossk" | "Lando Calrissian" |
| [25] "Lobot" | "Ackbar" | "Mon Mothma" |
| [28] "Arvel Crynyd" | "Wicket Systri Warrick" | "Nien Nunb" |
| [31] "Qui-Gon Jinn" | "Nute Gunray" | "Finn Valorum" |
| [34] "Jar Jar Binks" | "Roos Tarpals" | "Rugor Nass" |
| [37] "Ric Olié" | "Watto" | "Sebulba" |
| [40] "Quarsh Panaka" | "Shmi Skywalker" | "Darth Maul" |
| [43] "Bib Fortuna" | "Ayla Secura" | "Dud Bolt" |
| [46] "Galgano" | "Ben Quadinaros" | "Mace Windu" |
| [49] "Ki-Adi-Mundi" | "Kitt Fisto" | "Eeth Koth" |
| [52] "Adi Gallia" | "Saesee Tiin" | "Yarael Poof" |
| [55] "Plo Koon" | "Mas Amedda" | "Gregar Typho" |
| [58] "Cordé" | "Cliegg Lars" | "Poggle the Lesser" |

Exploring our dataset

Let's get to grips with our dataset.

How many animals do we have for each diet?

We can calculate `count` if we sum the

```
1 count(msl_eep, vore)
```

```
# A tibble: 5 × 2
  vore    <chr>    <int>
1 carni      19
2 herbi     32
3 insecti    5
4 omi       20
5 <NA>       7
```

count documentation

If we consult the documentation we find interesting

%>%

```
1 starwars %>%  
2 count(species, homeworld, sort = TRUE)
```

The code we will use to count the

```
1 bechdel %>%  
2 count(binary)
```

Let's address what ...

Little Bunny Foo Foo

To introduce pipes, we're going to **highlight** a **Wikipedia** article
How can we convert this poem into code?

```
Little bunny Foo Foo
Went hopping through the forest
Scooping up the
And bopping them on the head
```

Coding up little bunny

Little bunny Foo Foo
Went hopping through the forest
Scooping up the
And bopping them on the head

Let's create an instance of a bunny called

```
1 foo_foo <- little_bunny()
```

Now let's write the poem out as code:

```
1 bop_on(  
2   scoop_up(  
3     hop_through(foo_foo, forest),  
4     field_mouse  
5   ),  
6   head  
7 )
```

Understanding our code

In order to understand what our code does,

- Find the deepest `set` `path` `that happens`
- Work backwards (or up) the code

```
1 bop_on(  
2   scoop_up(  
3     hop_through(foo_foo, forest),  
4     field_mouse  
5   ),  
6   head  
7 )
```

This is exactly counter to the order of operation

Pi ping little bunny

Let's instantiate a `foo_bunny` called

```
1 foo_foo <- little_bunny()
```

Now write the same code as before but using

```
1 foo_foo %>%  
2   hop_through(forest) %>%  
3   scoop_up(field_mouse) %>%  
4   bop_on(head)
```

The order we read operations is exactly the happen!

Comparing the two

Independent of pipes, we create ourselves a

```
1 foo_foo <- little_bunny()
```

Now comparing the two code samples, the one

```
1 bop_on(  
2   scoop_up(  
3     hop_through(foo_foo, forest),  
4     field_mouse  
5   ),  
6   head  
7 )
```

```
1 foo_foo %>%  
2   hop_through(forest) %>%  
3   scoop_up(field_mouse) %>%  
4   bop_on(head)
```

Admittedly, it doesn't actually do anything!

Simpler %>% example

The pipe operator takes the left-hand side argument of the right-hand side of the expression

```
1 "cats" %>% rep(4)
```

Here's another example of what is doing the work
%>% is an example of it's ungarbage code easier

Pushing the pipe further

In some cases you don't want to use a pipe further

```
1 "cats" %>% paste(".", "are great", "but one can have too many", ".")
```

Where does some from?

magrit is the package, %t% is a function in 2 has become ridiculous popular

The pipe is now a tidyverse idiom, it's a standard way of writing code

If you want to use %>% in your own scripts, you can use the following

```
1 use_pipe() %>%
```

Advice on using

The pipe isn't a hammer to be used without and read with pipes.

Try to break pipe chains into blocks of similar understanding at a glance:

```
1 raw_data <- read_csv("data-raw/the-file.csv")
2
3 clean_data <- raw_data %>%
4   clean() %>%
5   clean_it() %>%
6   cleaning() %>%
7   cleaned()
8
9 clean_data <- clean_data %>%
10  normalise() %>%
11  normaler() %>%
12  norm_it()
```

What if I? hate

I st' perfectly accept able to hate

Th æ fin'e .

I st' just sugar to sweeten the already lovely

Howe,ye,you need a basic understanding of it
ti dy ver(saen d beyond) .

Hierarchical counting

We can count by as many attributes as we like

```
1 msleep %>%  
2 count(order, vore, sort = TRUE)
```

```
# A tibble: 32 × 3  
  order      vore      n  
  <chr>    <chr>    <int>  
1 Rodentia herbi    16  
2 Carnivora carni    12  
3 Primates omi     10  
4 Artiodactyla herbi    5  
5 Cetacea carni    3  
6 Perissodactyla herbi    3  
7 Rodentia <NA>    3  
8 Soricomorpha omi    3  
9 Chiroptera insecti  2  
10 Hyracoidea herbi    2  
# ... with 22 more rows
```


Is there any missing

The `filter()` function allows us to query a dataset

```
1 msleep %>%  
2 filter(sleep_total > 12)
```

We use for equivalence tests:

```
1 msleep %>%  
2 filter(vore == "carni")
```

We can negate conditions in two different ways:

```
1 msleep %>%  
2 filter(vore != "carni")
```

```
1 msleep %>%  
2 filter(!vore == "carni")
```

Is there any missing

We can't use an `isna()` test. `isna().sum()`

```
1 msleep %>%  
2 filter(is.na(conservation))
```

The `drop_na()` function returns only `NA` rows

```
1 msleep %>%  
2 drop_na()
```

{ n a n i a r } a n d { g g p l o t }

We're going to introduce { ggplot2 } today with scratch.

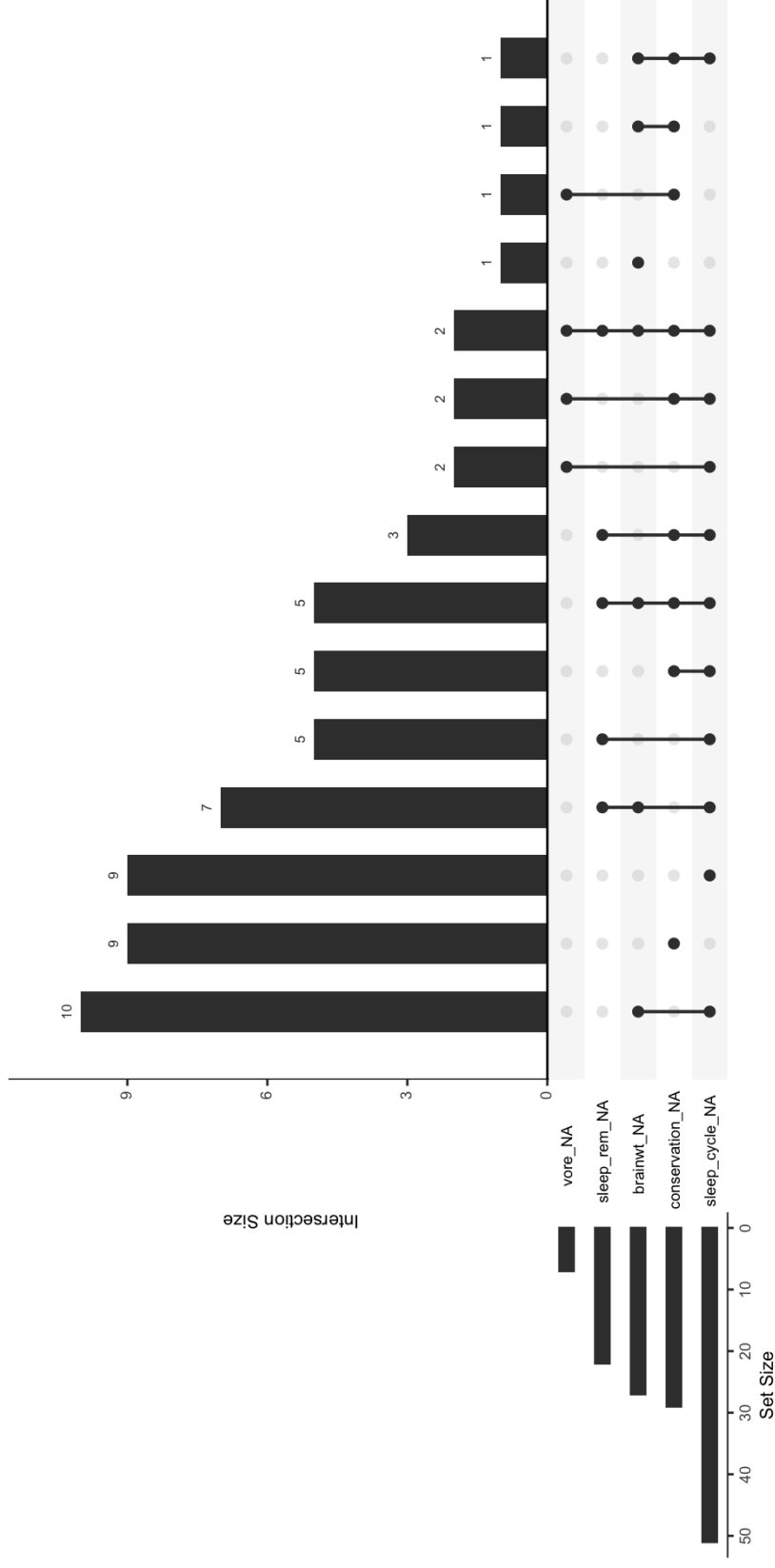
There are lots of effective abstractions of ggplot

Some of these packages provide these packages additional help you build a visualization charts. { naniar } is an example

{naniar} and {ggplot}

1. Install and load the {naniar} package
2. Run this code:

```
1 msleep %>%  
2 gg_miss_upset()
```



Making our own data

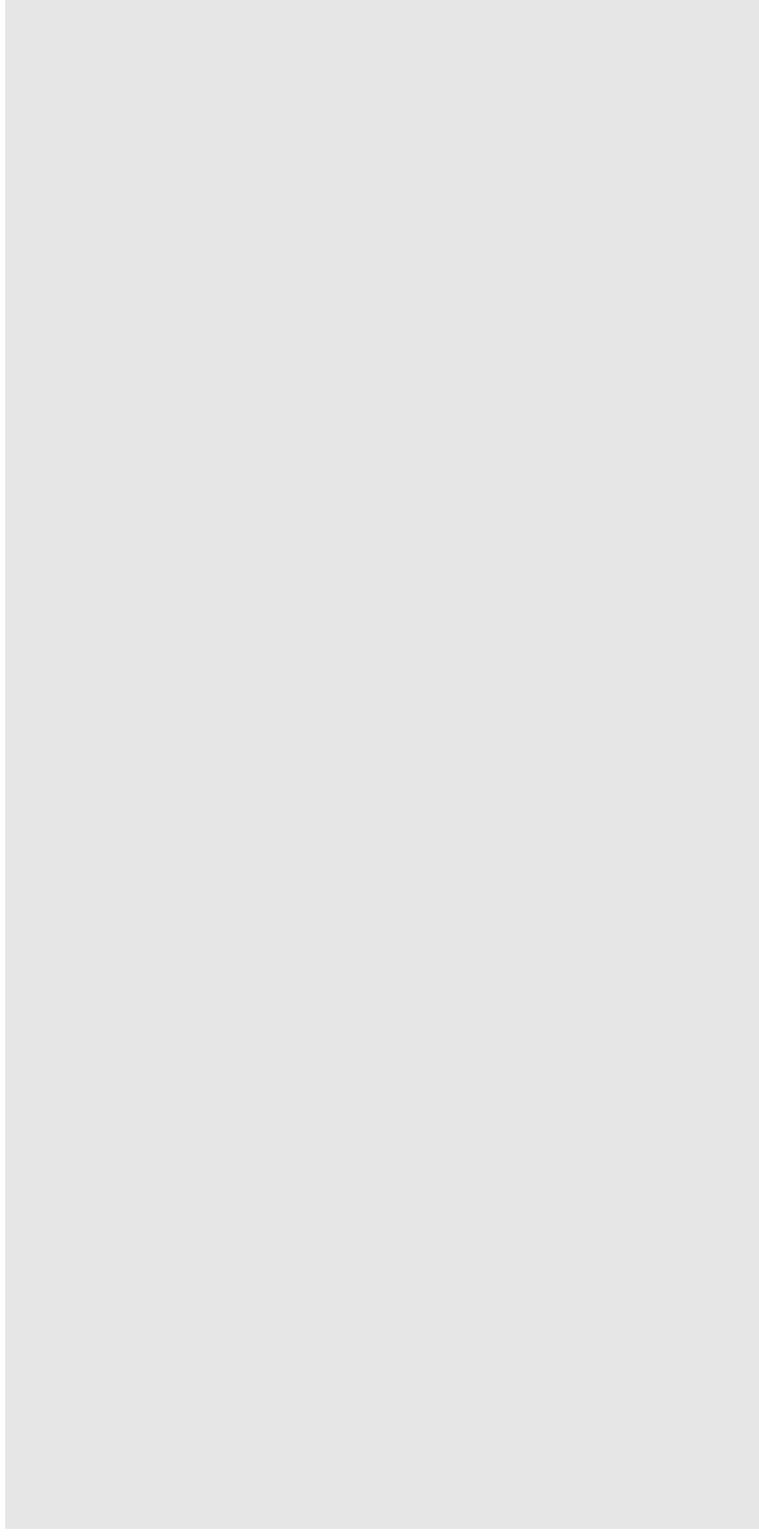
We're going to create a dataset with the following variables:

- Scatter plot of sleep_rem vs sleep_total
- Bar chart of mean sleep_total pervore

sleep_rem scatter plot

We start ggplot2 charts by providing a data

```
1 msl_eep %>%  
2 ggplot()
```



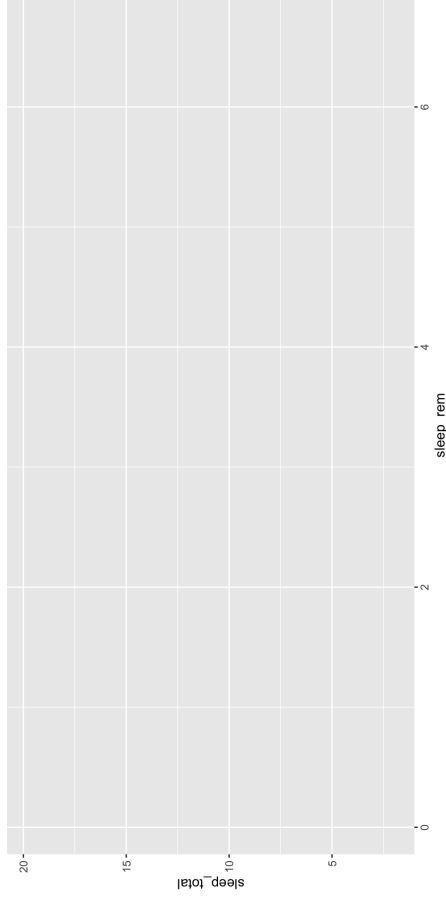
This creates an informative plot
so it can create a meaningful data viz.

sleep_rem scatter plot

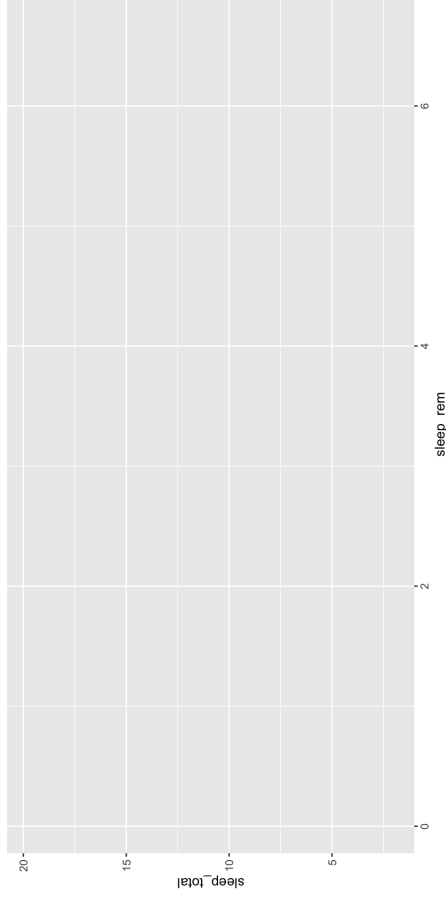
We need to tell {ggplot2} how to map column chart.

We do this with the aes() function. we have two different ways:

```
1 msleep %>%  
2   ggplot(aes(x = sleep_rem  
3             y = sleep_total))
```



```
1 msleep %>%  
2   ggplot() +  
3   aes(x = sleep_rem  
4       y = sleep_total)
```



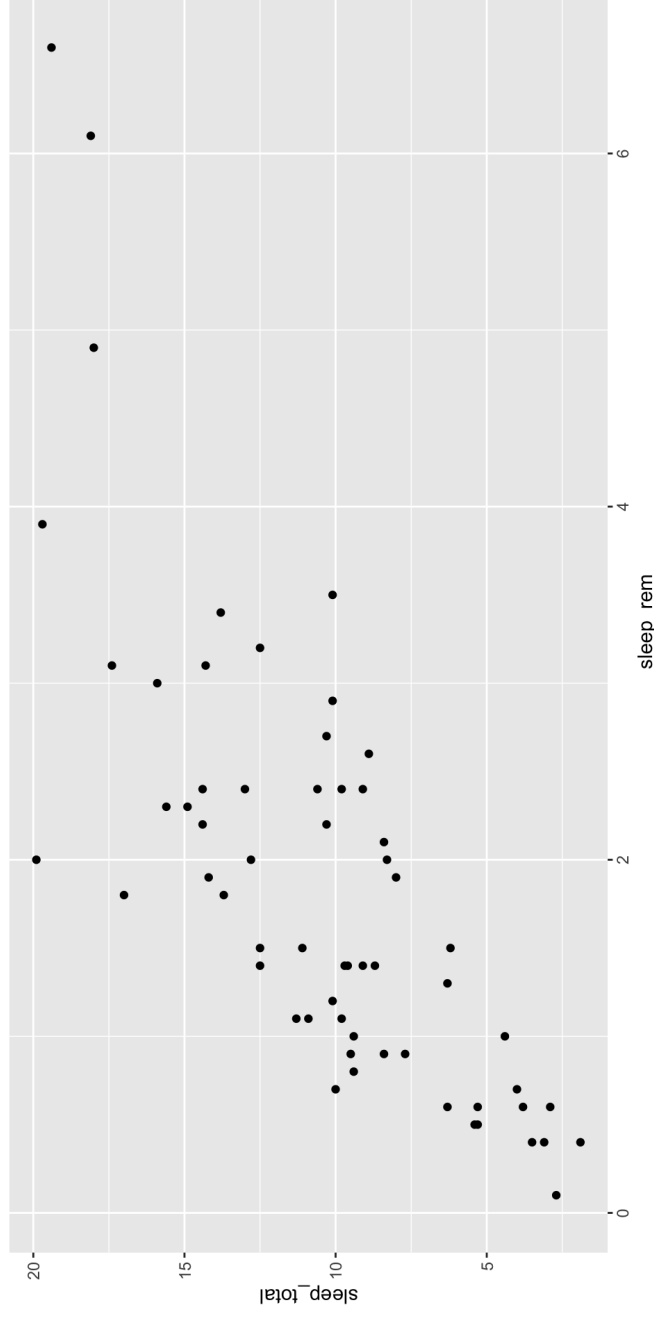
- { ggplot2 } uses the `scales` package to handle continuous scales
- The { tidyverse } functions are written specifically to evaluate

sleep_rem scatter plot

We now add geoms to our charts. These use shapes to our chart.

- { ggplot2 } ~~have~~ ~~retrieved~~ ~~added~~ ~~ways~~

```
1 msleep %>%  
2   ggplot() +  
3   aes(x = sleep_rem  
4       y = sleep_total) +  
5   geom_point()
```



sleep_rem_scatter plot

This is a pretty useless chart. It doesn't

What can we do to improve the chart?

sleep_reme_bar chart :

Before we can create a bar chart of the mean these values!

This means if you group by vore then calculate the mean

1. Add groups to data

```
1 msl_eep %>%  
2 group_by(vore)
```

```
# A tibble: 83 x 11  
# Groups:   vore [5]  
name      genus vore order conse...1 sleep...2 sleep...3 sleep...  
awake brainwt  
<chr> <chr> <chr> <chr> <chr> <dbl> <dbl> <dbl>  
<dbl> <dbl>  
1 Cheetah Aci n... carni Carn... l c 12.1 NA NA  
11.9 NA  
2 Owl monkey Aotus omi Pri m.. <NA> 17 1.8 NA 7  
0.0155  
3 Mountain be... Apl o... herbi Rode... nt 14.4 2.4 NA  
9.6 NA  
4 Greater sho... Bl ar ... omi Sori ... l c 14.9 2.3 0.133  
9.1 0.00029  
5 Cow Bos herbi Arti ... domest ... 4 0.7 0.667 20  
0.423  
6 Three-toed ... Brad... herbi Pil o... <NA> 14.4 2.2 0.767  
9.6 NA  
7 Northern fu... Cal l ... carni Carn... vu 8.7 1.4 0.383  
15.3 NA  
8 Vesper mouse Cal o... <NA> Rode... <NA> 7 NA NA 17  
NA  
9 Dog Cani s carni Carn... domest ... 10.1 2.9 0.333  
13.9 0.07
```

2. Calculate i

13 Ungroup the

finished.

sleep_remembar chart: sleep_remembar

Before we can create a bar chart of the mean values!

This means if you group by vore, you can calculate the mean sleep for each group

1. Add group

```
1 msleep %>%
2   group_by(vore) %>%
3   mutate(mean_sleep_total = mean(sleep_total))
```

2. Calculate in group

• mutate() leaves

3. Ungroup to find sd

```
# A tibble: 83 x 12
# Groups:   vore [5]
  name      genus vore order conse...1 sleep...2 sleep...3 sleep...
  <chr>      <chr> <chr> <chr> <chr> <dbl> <dbl> <dbl>
1 Cheetah    Acin... carni Carn...lc 12.1 NA NA
11.9 NA
2 Owl monkey Aotus omi Prim...<NA> 17 1.8 NA 7
0.0155
3 Mountain be... Apl o... herbi Rode...nt 14.4 2.4 NA
9.6 NA
4 Greater sho... Blar... omi Sori...lc 14.9 2.3 0.133
9.1 0.00029
5 Cow        Bos herbi Arti...domest... 4 0.7 0.667 20
0.423
6 Three-toed ... Brad... herbi Pil o... <NA> 14.4 2.2 0.767
9.6 NA
7 Northern fu... Cal l ... carni Carn...vu 8.7 1.4 0.383
15.3 NA
8 Vesper mouse Cal o... <NA> Rode... <NA> 7 NA NA 17
NA
9 Dog        Canis carni Carn...domest... 10.1 2.9 0.333
13.9 0.07
```

sleep_remem_bar_chart:

Before we can create a bar chart of the mean values!

This means `group_by(vore)` for calculation

```
1 msleep %>%  
2   group_by(vore) %>%  
3   summarise(mean_sleep_total = mean(sleep_total))
```

```
2. Calculate in group  
# A tibble: 5 x 2  
  vore mean_sleep_total  
  <chr>          <dbl>  
1 carni          10.4  
2 herbi           9.51  
3 insecti        14.9  
4 omni           10.9  
5 <NA>           10.2
```

- `mutate()` leaves columns
- `group_by()` throws away columns

3. Ungroup `final_sd`

sleep_reme_bar chart :

Before we can create a bar chart of the mean these values!

This means if you group by vore then calculate the mean sleep

1. Add groups to

2. Calculate i

• mutate() leave

columns

• group_by() then
columns

3. Ungroup then summarise

```
1 msleep %>%
2   group_by(vore) %>%
3   summarise(mean_sleep_total = mean(sleep_total)) %>%
4   ungroup()
```

```
# A tibble: 5 x 2
  vore mean_sleep_total
<chr> <dbl>
1 carni 10.4
2 herbi 9.51
3 insecti 14.9
4 omi 10.9
5 <NA> 10.2
```

sleep_reme_bar chart :

Before we can create a bar chart of the mean values!

This means if you want to calculate the mean

1. Add groups to

```
1 msl_eep %>%
2   group_by(vore) %>%
3   summarise(mean_sleep_total = mean(sleep_total)) %>%
4   ungroup()
```

2. Calculate i

```
# A tibble: 5 x 2
  vore mean_sleep_total
<chr> <dbl>
1 carni 10.4
2 herbi 9.51
3 insecti 14.9
4 omi 10.9
5 <NA> 10.2
```

• mutate() leave

columns

• group_by() th

columns

3. Ungroup the data when

Don't forget to m

assignment for th

subset/analysis of the dataset

```
1 mean_sleep_by_vore <- msl_eep %>%
2   group_by(vore) %>%
3   summarise(mean_sleep_total = mean(sleep_total)) %>%
4   ungroup()
```

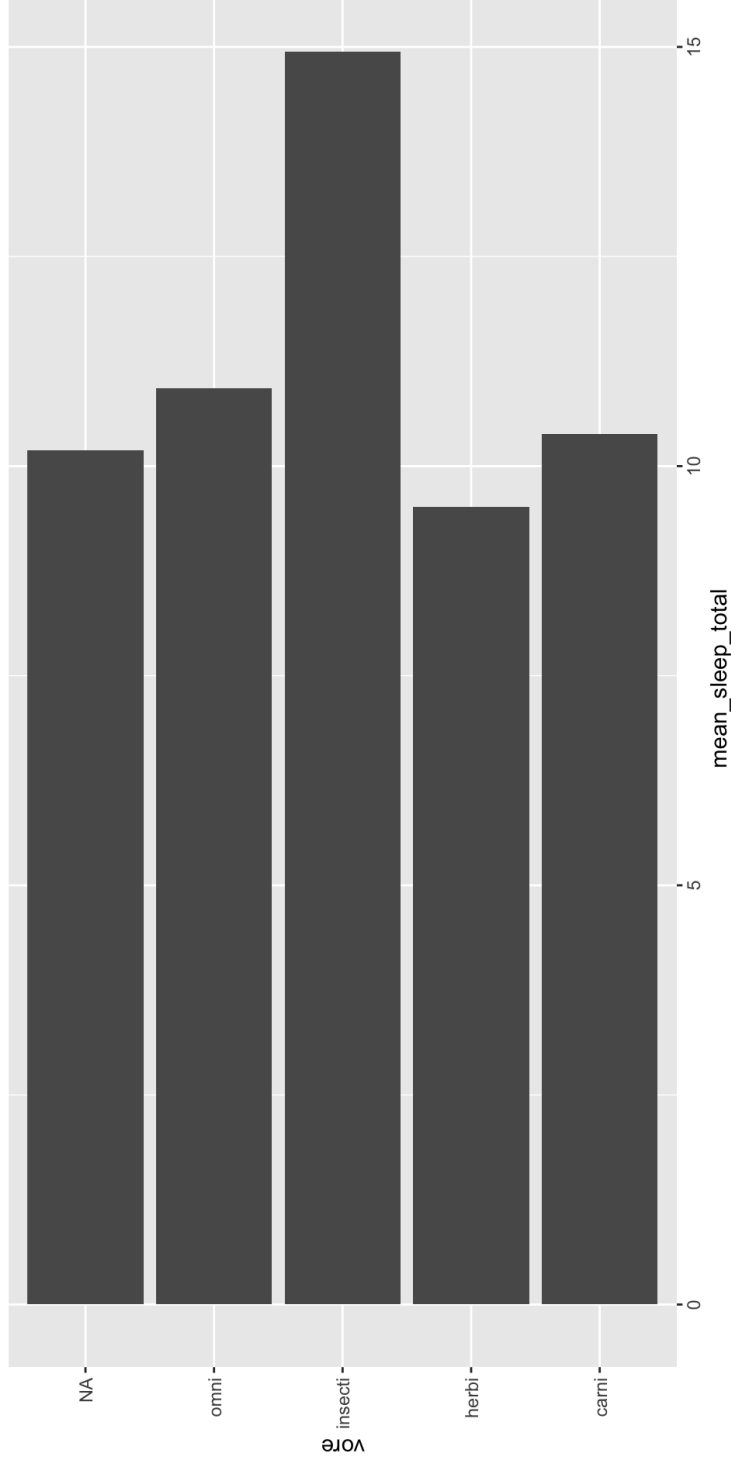
sleep_rem_bar_chart:

Please setup a ggplot with the following

- x axis should be the "mean_sleep_total" col
- y axis should be the "vore" column

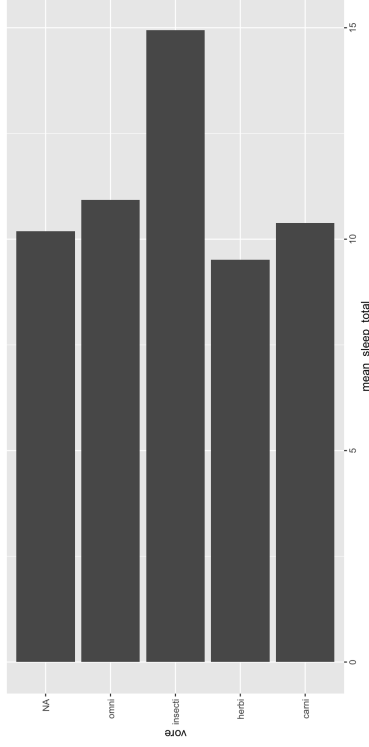
sleep_remap_bar chart:

```
1 mean_sleep_by_vore %>%  
2   ggplot() +  
3   aes(x = mean_sleep_total,  
4       y = vore) +  
5   geom_col()
```



How can we improve this chart?!

sleep_rem bar chart :



This chart is confusing to bars are shown.

{ ggplot2 } uses information the order of scales.

- In the `character` column ordering is used
- In the `factor` columns the levels in the factor

We're going to come back visualisation week.

Mutating multiple columns

Often we need to target multiple columns at once. It's a bit tedious to write the code to fractionally

This is achieved with the `mutate()` function.

The `across()` function takes a list of column names or a vector of column indices and applies the same function to all of them.

```
1 msleep %>%  
2 mutate(across(argument_1, argument_2))
```

Mutating multiple columns

Often we need to target multiple columns at once. It's a `pd.DataFrame` in the dataset to fraction

This is achieved with `with` but it's

The first argument is `fraction`

This is achieved with `fraction`.

```
1 msleep %>%  
2 mutate(across(argument_1, argument_2))
```

```
1 msleep %>%  
2 mutate(across(starts_with("sleep"),  
3 argument_2))
```

Multiplying multiple columns

Often we need to target multiple columns at once. It's a pain to write the code to fractionally

This is achieved with the

The 2nd argument is where to each column header:

We need to write a function here.

```
1 msleep %>%
2 mutate(across(argument_1, argument_2))
```

Thus usualcyhieved with that and

```
1 ~. / 24
```

```
1 msleep %>%
2 mutate(across(starts_with("sleep"),
3 ~. / 24))
```

```
# A tibble: 83 x 11
  name      genus vore order conse...!
sleep...? sleep...3 sleep... awake brainwt
<chr> <chr> <chr> <chr> <chr> <chr>
<dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 Cheetah  Acin... carni Carn...lc
0.504 NA NA 11.9 NA
2 Owl monkey Aotus omi Prim...<NA>
0.708 0.075 NA 7 0.0155
3 Mountain b... Apl o... herbi Rode... nt
0.1 NA 9.6 NA
4 Greater sh... Bl ar... omi Sori...lc
0.621 0.0958 0.00556 9.1 0.00029
5 Cow Bos herbi Arti... domest...
0.167 0.0292 0.0278 20 0.423
6 Three-toed... Brad... herbi Pil o... <NA>
0.0917 0.0319 9.6 NA
7 Northern f... Cal l... carni Carn... vu
0.362 0.0583 0.0160 15.3 NA
8 Vesper mou... Cal o... <NA> Rode... <NA>
0.292 NA NA 17 NA
9 Dog Canis carni Carn... domest...
0.421 0.121 0.0139 13.9 0.07
10 Roe deer Canr herbi Arti lc
```

Mutating multiple columns

Often we need to target multiple columns at once. It's a pain in the dataset to fraction

This is achieved with `mutate()` but it's

The 2nd argument is where
to each column head:

We need to write a function here.

Thus usually we don't need
to write a function here.

1 ~ . / 24

All tidyverse wrangling is done "column-wise" operations. [visit the tidyverse website](#) to see some examples of it.

g a p m i n d e r



Task: Get the gapminder

SLIDE 1 OF 3

1. Add a heading for the gapminder dataset.
2. Load the `gapminder` package in the setup code chunk
3. Add a new code chunk `gapminder` to read in the data

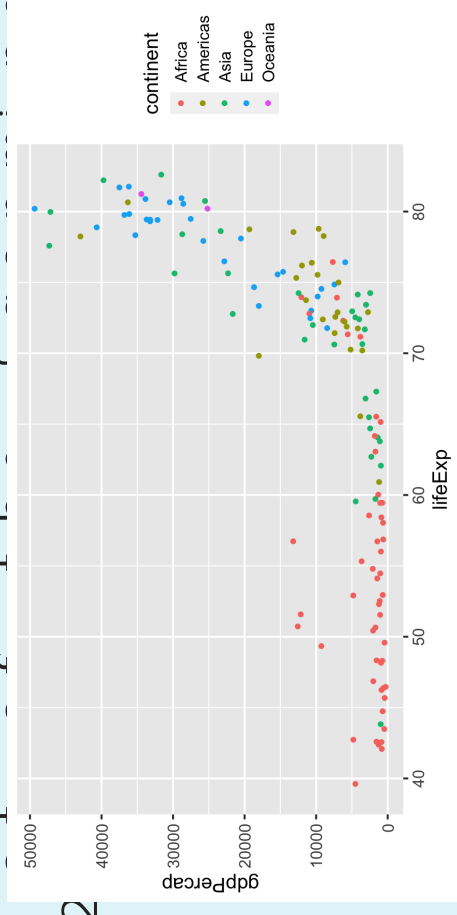
1 gapminder



task: gapminder scatter

SLIDE 2 OF 3

Create this scatter plot
dataset for the year 2007

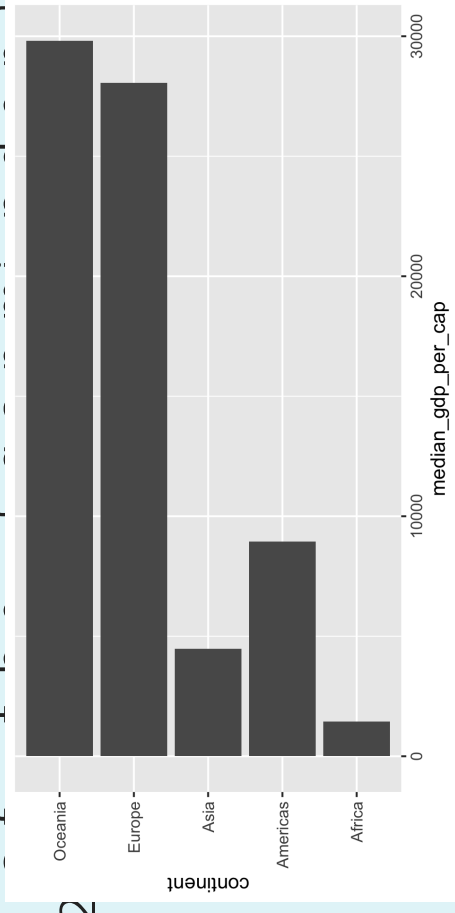




task: gapminder bar chart

SLIDE 3 OF 3

Create this bar chart
dataset for the year 2007



GBD Dataset



Task: Get the GBD data

SLIDE 1 OF 2

1. Add a sub-folder `data` your project called
2. Inside folder add `oabtsadni-pdata` called
3. Add this code

```
1 download.file("https://raw.githubusercontent.com/charliejhadley/eng7218_data-science-for-healthcare-applications_bcu-  
2 destfile = "data/global-burden-of-disease-data.csv")
```

5. Run the code



Task: Get the GBD data

SLIDE 2 OF 2

1. Add a new heading for the GBD Dataset to y

Reading data into R

The `readr` package provides excellent `read_csv()` and `read_delimited_csv()` files. I discuss `read_csv()`.

We need to think about the `read_csv()` function. The following:

1. Add a code chunk to your
2. Choose a name for the data frame you create to store something similar to denote this is your

```
1 di sease_burden_raw <-
```

3. Call the appropriate function to read the data

```
1 di sease_burden_raw <- read_csv("")
```

4. Print the cursor inside the `read_csv()` and select your

```
1 di sease_burden_raw <- read_csv("data/global-burden-of-di sease-data.csv")
```

Matching text in R (

A lot of data wrangling, most of the computer science we usually fear isn't exactly "user-friendly" character".

The GBD dataset gives
example of this.

There are not different
the dataset

- World Bank regions
- Geographic regions

The tidyverse gives us `satr` in `en` to
strings.

```
1 disease_burden_raw %>%  
2 count(location_name)
```

```
# A tibble: 10 x 2  
  location_name      n  
  <chr>          <int>  
1 African Region    720  
2 Eastern Mediterranean Region 720  
3 European Region   720  
4 Region of the Americas        720  
5 South-East Asia Region        720  
6 Western Pacific Region        720  
7 World Bank High Income        720  
8 World Bank Low Income         720  
9 World Bank Lower Middle Income 720  
10 World Bank Upper Middle Income 720
```


Matching text in R (

We can search the beginning of strings with `str_starts()`

```
1 di sease_burden_raw %>%  
2 filter(str_starts(location_name, "World Bank"))
```

```
# A tibble: 2,880 × 16  
  measu...1 measu...2 locat...3 sex_id sex_n... age_id  
age_n... cause... cause...  
  <chr> <dbl> <chr> <dbl> <chr> <dbl> <chr> <dbl>  
1 1 Deaths 44575 World ... 3 Both 22  
All ag... 409 Non-co...  
2 1 Deaths 44575 World ... 3 Both 22  
All ag... 409 Non-co...  
3 1 Deaths 44575 World ... 3 Both 22  
All ag... 409 Non-co...  
4 1 Deaths 44575 World ... 3 Both 22  
All ag... 294 All ca...  
5 1 Deaths 44575 World ... 3 Both 22  
All ag... 294 All ca...  
6 1 Deaths 44575 World ... 3 Both 22  
All ag... 294 All ca...  
7 1 Deaths 44575 World ... 3 Both 22  
All ag... 295 Commun...  
8 1 Deaths 44575 World ... 3 Both 22  
All ag... 295 Commun...  
9 1 Deaths 44575 World ... 3 Both 22  
All ag... 295 Commun...  
10 1 Deaths 44576 World ... 3 Both 22
```

Matching text in R (

For more complex string
have to make use of R
REGEX stands for regular
and is an approach to
that is implemented in
languages.

```
1 di sease_burden_raw %>%  
2 filter(str_detect(location_name, "^World Bank"))
```

```
# A tibble: 2,880 x 16  
  measu...1 measu...2 locat...3 locat... sex_id sex_n... age_id  
age_n... cause... cause...  
  <chr> <dbl> <chr> <dbl> <chr> <dbl> <chr> <dbl>  
1 1 Deaths 44575 World ... 3 Both 22  
All ag... 409 Non-co...  
2 1 Deaths 44575 World ... 3 Both 22  
All ag... 409 Non-co...  
3 1 Deaths 44575 World ... 3 Both 22  
All ag... 409 Non-co...  
4 1 Deaths 44575 World ... 3 Both 22  
All ag... 294 All ca...  
5 1 Deaths 44575 World ... 3 Both 22  
All ag... 294 All ca...  
6 1 Deaths 44575 World ... 3 Both 22  
All ag... 294 All ca...  
7 1 Deaths 44575 World ... 3 Both 22  
All ag... 295 Commun...  
8 1 Deaths 44575 World ... 3 Both 22  
All ag... 295 Commun...  
9 1 Deaths 44575 World ... 3 Both 22  
All ag... 295 Commun...  
10 1 Deaths 44575 World ... 3 Both 22
```

We can use regular exp
pattern{ts tfor n garl}l
functions.

regex101.scam really use
building up complex regular expressions.

GBD data visualization

I'd like to get the data so that all of these

- We can only see data for the most recent year
- We can see the death rate for each cause_name
- There are 161 countries in the dataset

GBD data visualisation

Once we download the data on the interesting

```
1 di sease_burden_percent_deaths %>%  
2 select(location_name, cause_name, val)
```

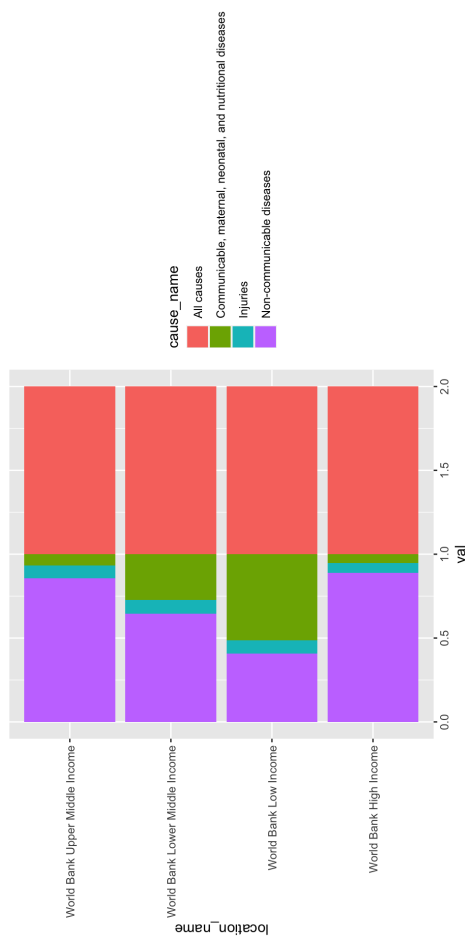
```
# A tibble: 16 x 3  
  location_name cause_name val  
  <chr> <chr> <dbl>  
1 World Bank High Income All causes 1  
2 World Bank High Income Communicable, maternal, neonatal, and ... 0.0528  
3 World Bank Upper Middle Income Non-communicable diseases 0.856  
4 World Bank Upper Middle Income All causes 1  
5 World Bank Upper Middle Income Communicable, maternal, neonatal, and ... 0.0662  
6 World Bank Low Income All causes 1  
7 World Bank Upper Middle Income Injuries 0.0780  
8 World Bank Low Income Communicable, maternal, neonatal, and ... 0.513  
9 World Bank High Income Injuries 0.0579  
10 World Bank Lower Middle Income Non-communicable diseases 0.645  
11 World Bank High Income Non-communicable diseases 0.889  
12 World Bank Low Income Injuries 0.0812  
13 World Bank Lower Middle Income All causes 1  
14 World Bank Lower Middle Income Communicable, maternal, neonatal, and ... 0.273  
15 World Bank Lower Middle Income Injuries 0.0824  
16 World Bank Low Income Non-communicable diseases 0.406
```

Let's visualise this as a bar chart

GBD data visualization statistics

How can we make this chart more meaningful

```
1 disease_burden_percent_deaths %>%
2 select(location_name, cause_name, val) %>%
3 ggplot() +
4 aes(x = val,
5      y = location_name,
6      fill = cause_name) +
7 geom_col()
```



Going further

We've only barely scratched the surface of the **data science** field. There are a lot of topics we haven't covered more advanced topics like **machine learning** and **deep learning**.

- Read in file `twi_atdh_*()` functions
- Understand `read_csv()` and `read_csv_from_text_file()`
- Filter `data.frame` with `filter()`
- Use `select()` search/modify columns
- Use `mutate()` modify existing columns and add new columns
- Use `group_by()` calculate in-group values
- Use `ggplot2` for exploratory data analysis.

References

1. Wickham, H. & Groland, G. R. for data science: Import, tidy transform, visualize, and model data. (O'Reilly, 2016).
2. Matloff, N. Teaching R in a kinder, gentler, more effective manner: (2022).
3. Global Health Data Exchange. Global Burden of Disease Dataset Explorer. Institute for Health Metrics and Evaluation (2022).
4. Savage, V. M. & West, G. B. **A quantitative and theoretical framework for understanding the scaling of the** National Academy of Sciences 104, 1051–1056 (2007).
5. Hans Rosling. **The best stats you've ever seen** [Video]. The best stats you've ever seen (2006).