

# Week 4: Data visualisation

Charlotte Hadley

# Topics for today

1. Why do we use charts to tell stories?
2. Evidence-based visual perception theory
3. Advice on choosing charts
4. Advice on using colour in charts
5. Using this advice to tell stories with charts built with `{ggplot2}`

# Why do we use charts??

A picture is worth a thousand words

---

# Data visualisations are demonstrably useful

There is considerable experimental evidence for data visualisations improving:

- Comprehension of data
- Decision making accuracy and confidence

Evidence has been collected using eye-tracking, survey filling and interviews.

For a good overview of the available research see Eberhard 2021<sup>1</sup>.

---

Some of these studies consider tables to be a type of data visualisation.

I agree with this! Tables are often awesome choices for presenting data - let's talk more about this later today.

# Data visualisations are demonstrably useful

In 1973 Anscombe<sup>2</sup> published a paper designed to demonstrate...

Graphs are essential to good statistical analysis.

Each of the four data sets yields the same standard output from a typical regression program, namely

Number of observations ( $n$ ) = 11

Mean of the  $x$ 's ( $\bar{x}$ ) = 9.0

Mean of the  $y$ 's ( $\bar{y}$ ) = 7.5

Regression coefficient ( $b_1$ ) of  $y$  on  $x$  = 0.5

Equation of regression line:  $y = 3 + 0.5x$

Sum of squares of  $x - \bar{x}$  = 110.0

Regression sum of squares = 27.50 (1 d.f.)

Residual sum of squares of  $y$  = 13.75 (9 d.f.)

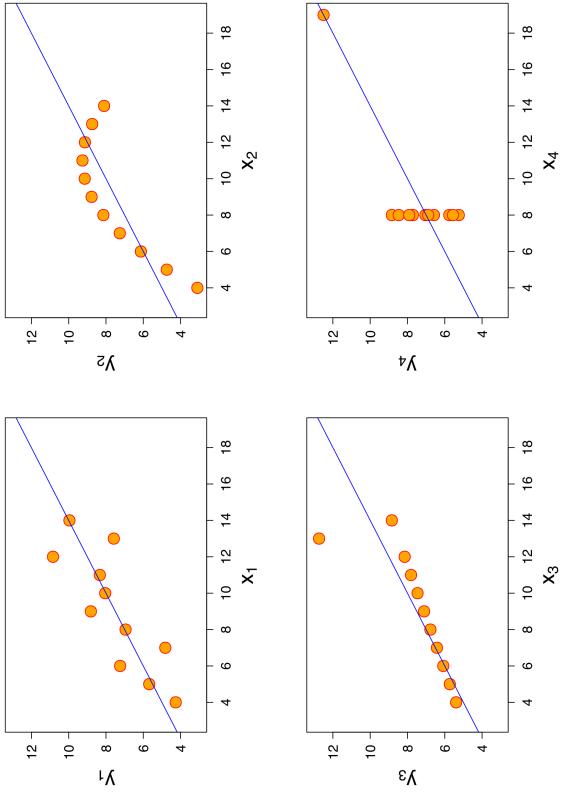
Estimated standard error of  $b_1$  = 0.118

Multiple  $R^2$  = 0.667

# Data visualisations are demonstrably useful

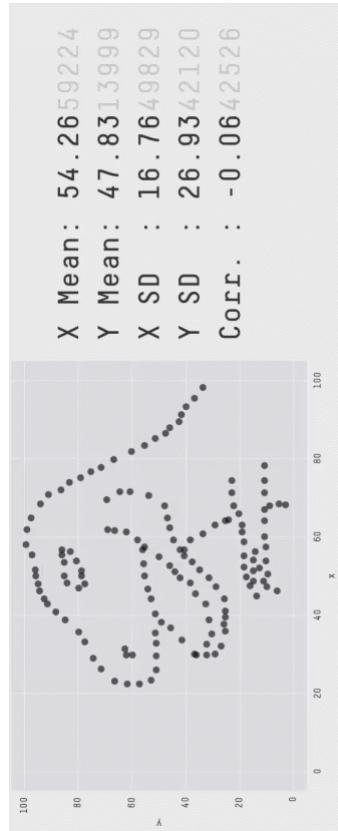
However, if you visualised the datasets it was obvious these datasets were fundamentally different to one another.

These charts are now known as *Anscombe's quartet*<sup>2</sup>.



# Data visualisations are demonstrably useful

The “Datasaurus Dozen” is a modern reimagining of the original quartet<sup>3</sup>.

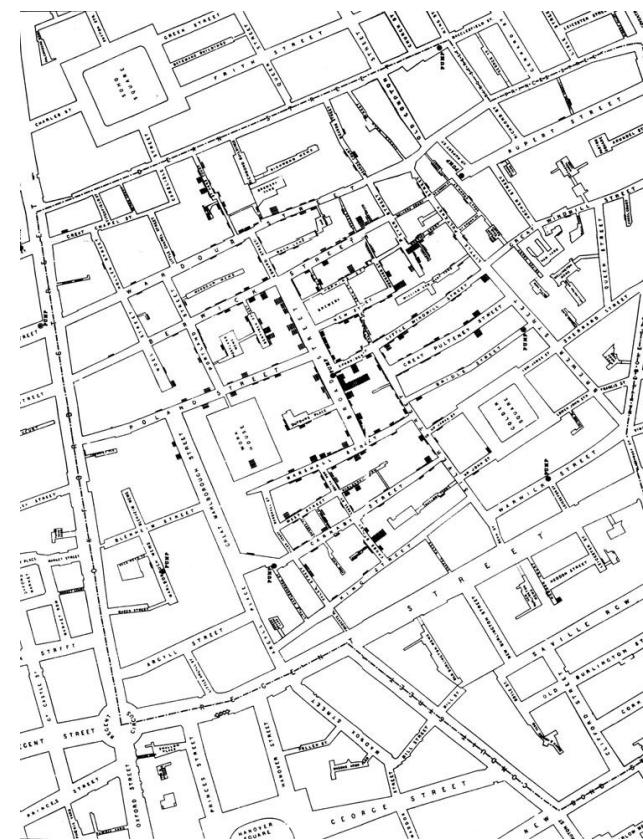


Datasaurus was originally created by Alberto Cairo<sup>4</sup>.

... there's now an R package for building your own metamers [eliocamp.github.io/metamer/](https://eliocamp.github.io/metamer/)

**Always visualise your datasets.**

# Data visualisations are demonstrably useful



There are several historical visualisations that have fundamentally changed social policy and behaviour.

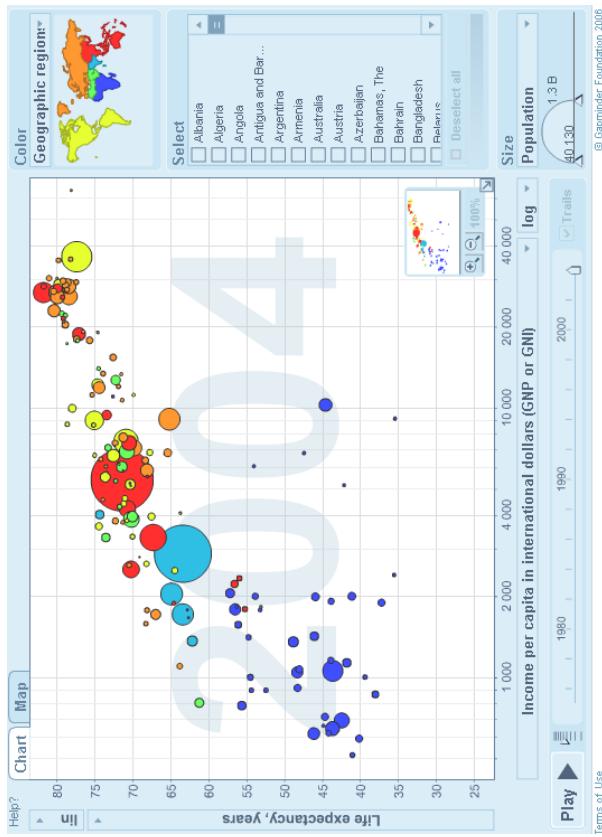
This is a map from John Snow in 1855<sup>5</sup> that ties a cholera outbreak to a specific water pump.

Combined with Snow's statistical analyses this was a significant step towards the development and acceptance of germ theory.

**Data visualisations are demonstrably useful**

# Data visualisations are demonstrably useful

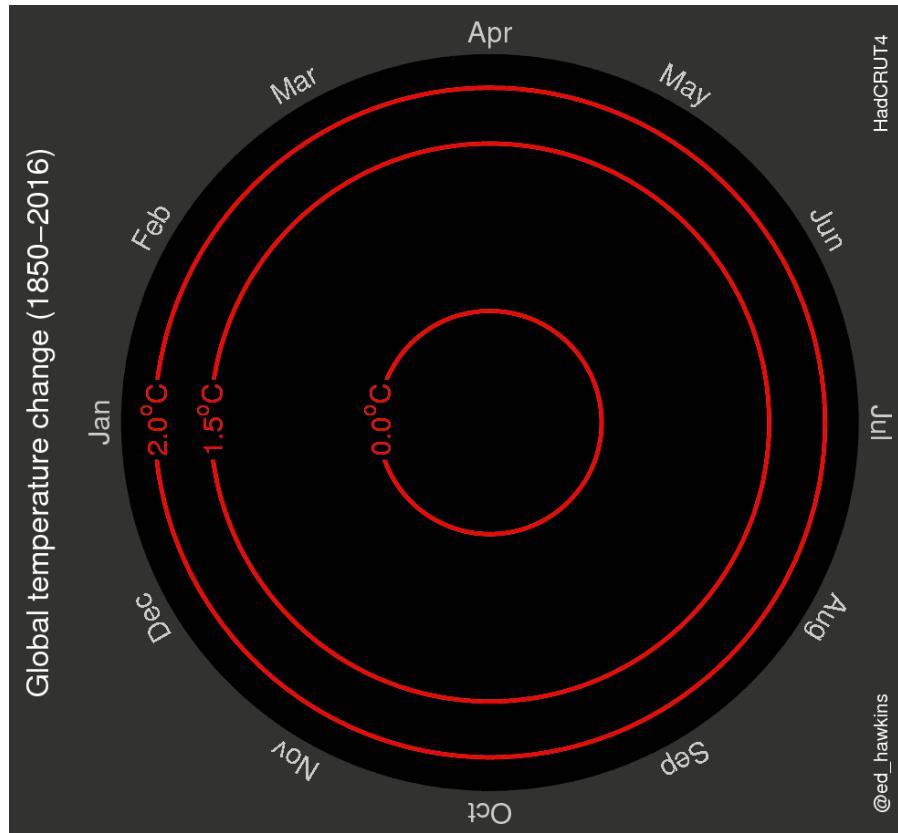
In 2006 Hans Rosling<sup>7</sup> gave an incredible **TED talk** where he introduced animated **bubble charts** as a tool to tell stories about global development.



These charts helped demonstrate the value of interactive and animated data visualisations - which is why Google bought the tool behind the charts!

# Data visualisations are demonstrably useful

A more recent example of a very powerful data visualisation is the spiralling global temperature GIF from 2016 by Ed Hawkins<sup>8</sup>.

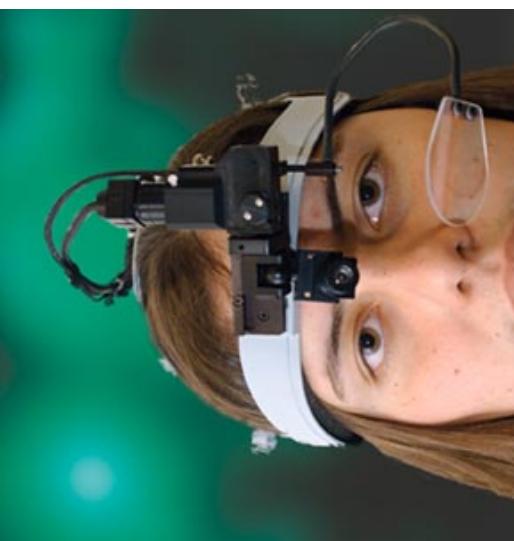


We can create animated GIF with `{ggplot2}` via the `{gganimate}` package. In fact, Pat Schloss<sup>9</sup> has a [YouTube video](#) and [GitHub repo](#) recreating this chart with R.

# Evidence-based visual perception theory

# Evidence-based visual perception theory

There is a wealth of evidence-based research in how precisely or accurately charts are perceived by readers.



Our evidence comes from:

- Eye tracking. We're really good at measuring where the eye is looking, for how long and how intently.
- Asking trial participants to estimate or compare values in charts.

Source: [Wikimedia.org](https://commons.wikimedia.org)

There are open debates<sup>1</sup> on how our internal visual perception system works - what the brain is doing.

# Elementary perceptual tasks

Back in 1984 Cleveland & McGill<sup>11</sup> published their seminal paper on graphical perception theory where they defined “elementary perceptual tasks”.

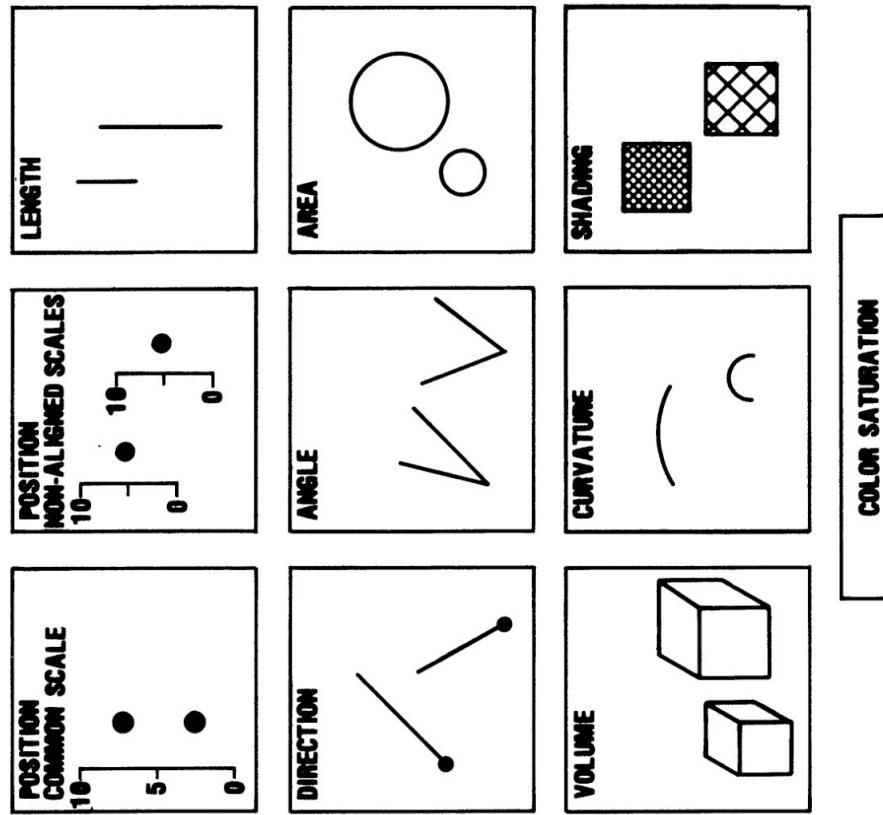


Figure 1. Elementary perceptual tasks.

# Elementary perceptual tasks

Cleveland & McGill<sup>11</sup> designed many experiments where participants were asked to:

- Identify the largest/smallest segment
- Estimate what % the smaller segment was of the larger segment

The accuracy of subject estimates was then statistically analysed.

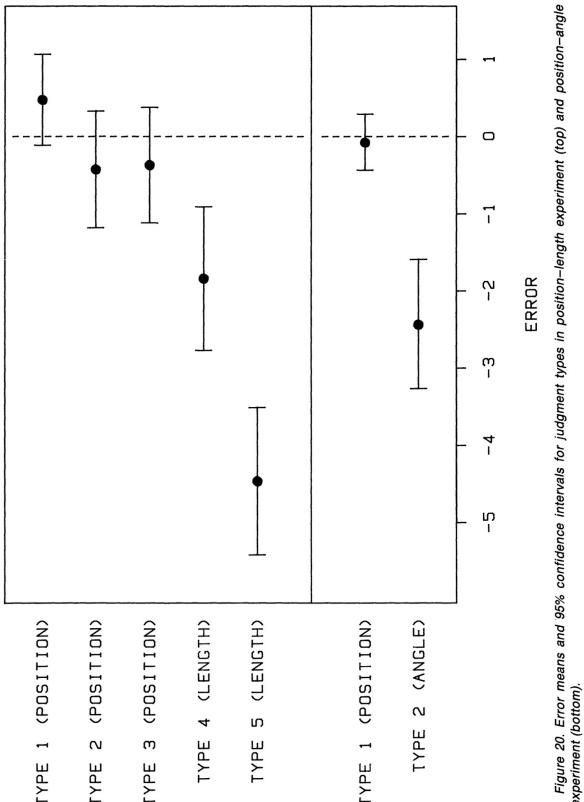
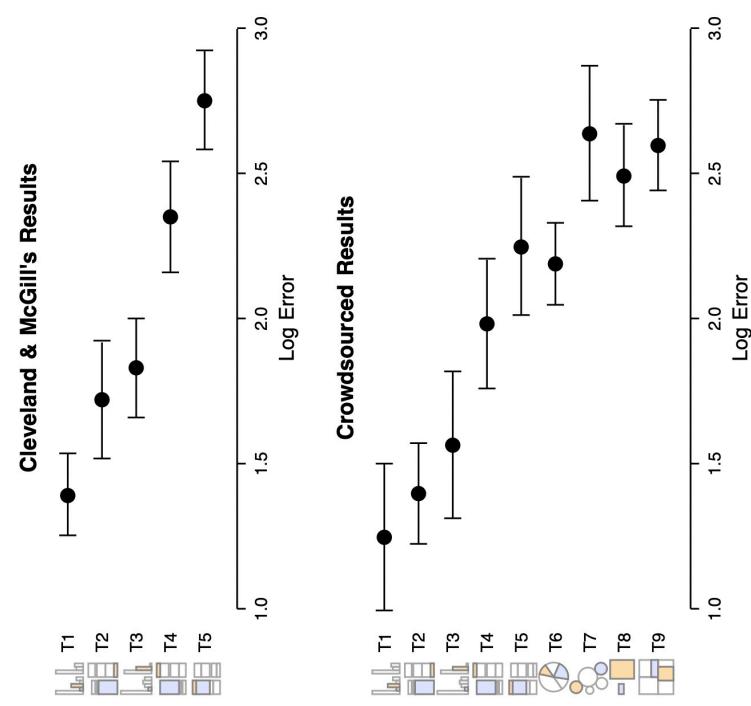


Figure 20. Error means and 95% confidence intervals for judgment types in position-length experiment (top) and position-angle experiment (bottom).

# Crowd-sourced evidence for perception theory

Heer & Bostock<sup>12</sup> replicated this study using Amazon's Mechanical Turk with 3,481 participants in 2010.

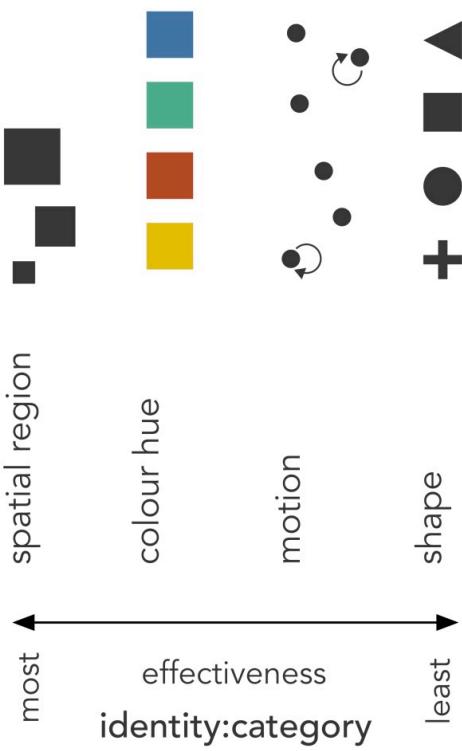
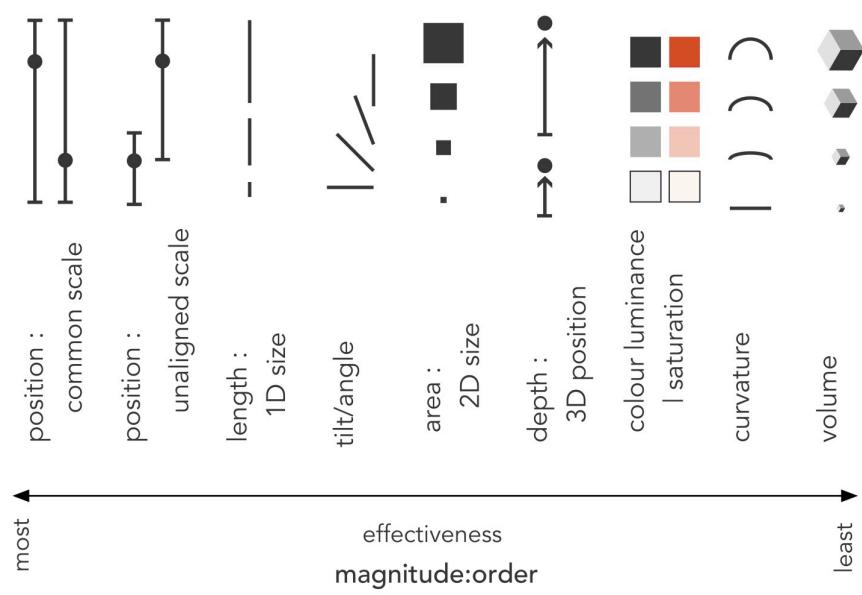


They validated the results of Cleveland & McGill<sup>11</sup> and provided further evidence that...

There is a hierarchy of elementary perceptual tasks - or chart elements - when accuracy matters.

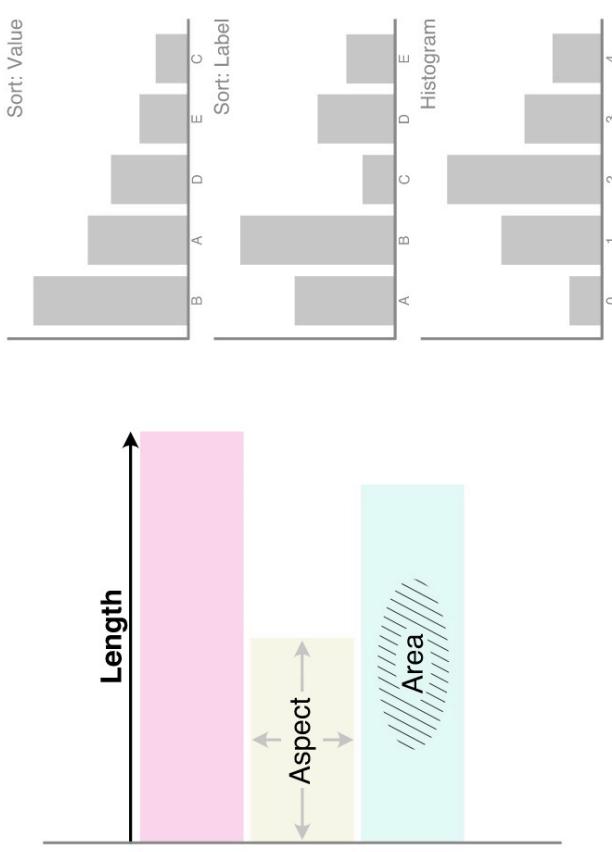
Figure 4: Proportional judgment results (Exp. 1A & B). Top: Cleveland & McGill's [7] lab study. Bottom: MTurk studies. Error bars indicate 95% confidence intervals.

# Ordering channels of communication (by accuracy)

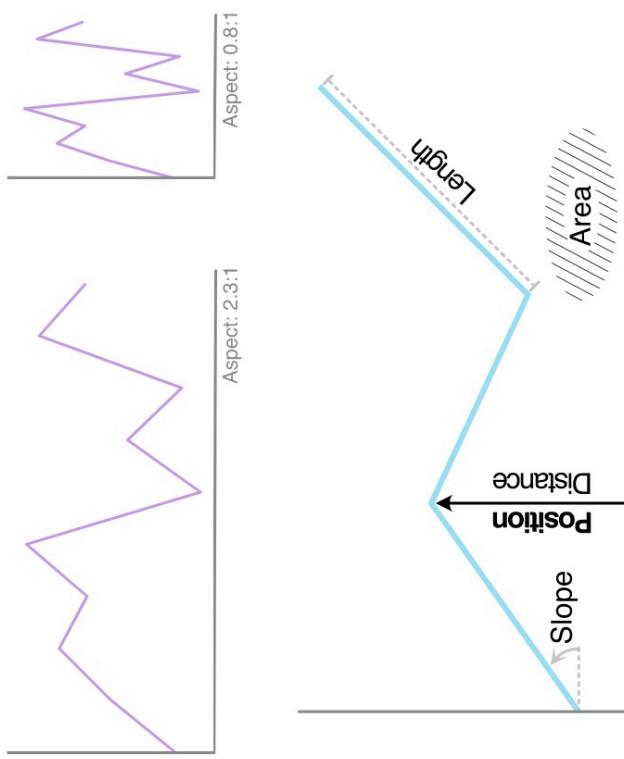


Images from Beecham et al<sup>13</sup>

## ... real-world applications of visual perception theory (I)



**FIGURE 2.** Bar charts encode values as their length, but also their area, aspect ratio, and overall shape. Sorting is often used for specific use cases.



**FIGURE 3.** Line charts are specified by the location of the points connected by lines, but are read as slope and length, as well as area. Aspect ratio of the chart is also generally considered important.

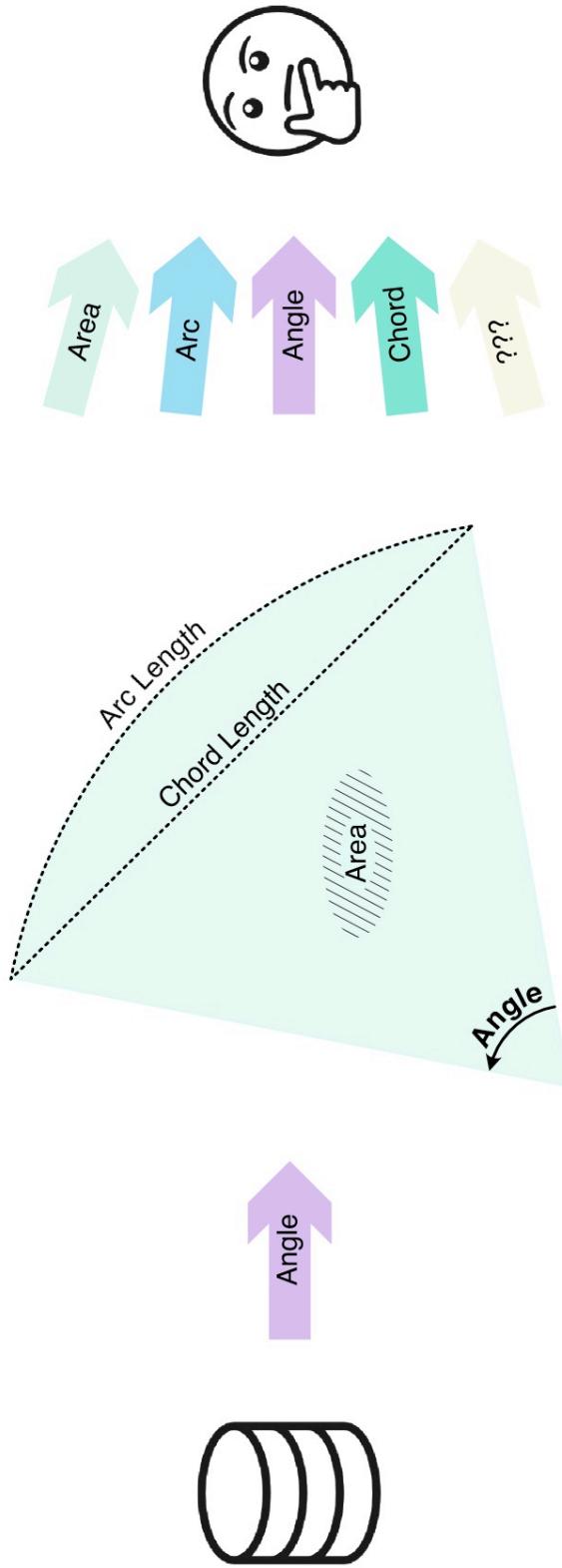
Images from Robert Kosara<sup>14</sup>

## ... real-world applications of visual perception theory (II)

The *specified encoding* transforms a data value into a shape, ...

... which expresses that value as multiple *observable encodings*, ...

... some subset of which end up being *observed* by the user.



**FIGURE 1.** Pie charts are specified by angle, but may also be read by area, arc length, or even chord. Shape recognition is also likely for specific angles like  $90^\circ/25\%$  and  $180^\circ/50\%$ .

Images from Robert Kosara<sup>14</sup>

## ... real-world applications of visual perception theory (III)

Data from USDA  
Graph by @irg\_bio

### Seasonality in the number of bee colonies

The number of bee colonies in the US changes with the seasons depending on the state. In some northern states, the number of bees increases more than two-fold in spring, and their numbers are maintained until winter. In some warmer southern states, such as California or Texas, the number of colonies drops during the hottest months, recovering only with milder temperatures.

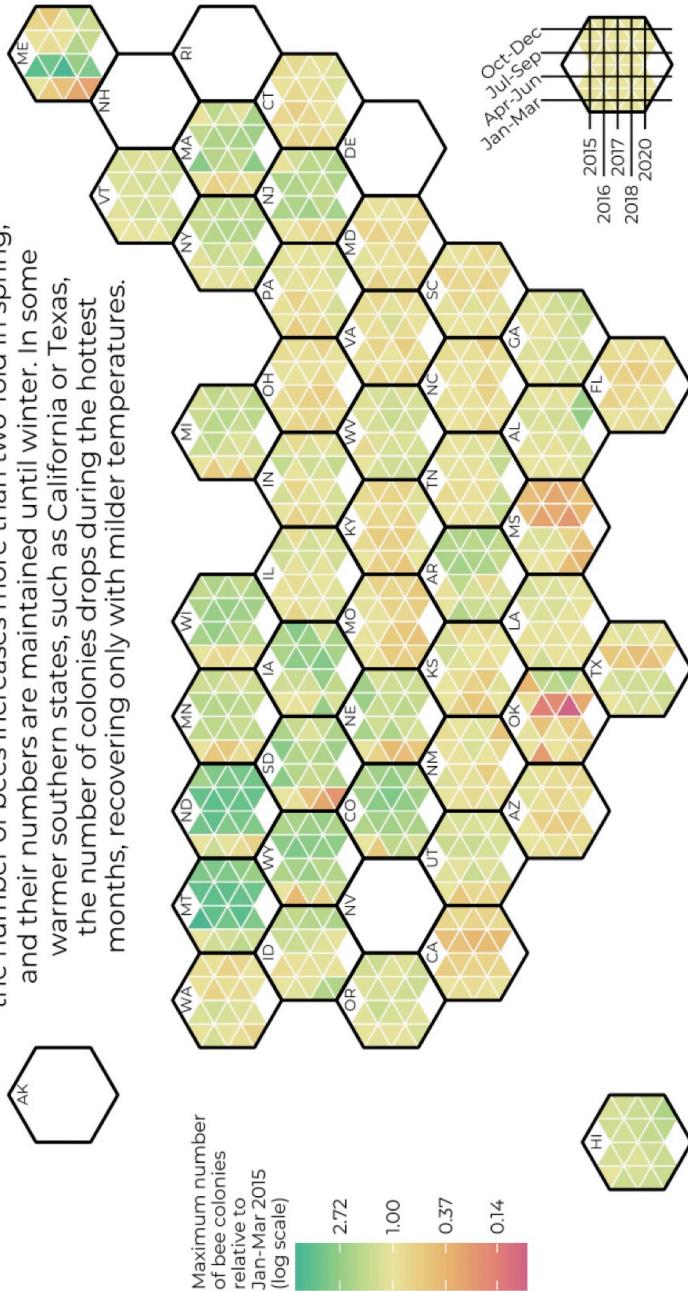
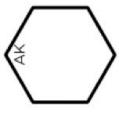


Image found on Twitter from [@irg\\_bio<sup>15</sup>](#) - code for chart available from [GitHub<sup>16</sup>](#).

Why is someone reading  
measuring your chart?

# Why is someone reading measuring your chart?

To extract accurate values

The magnitude of chart elements.

To quantitatively compare values.

The part to whole or relative magnitude of chart elements.

To find the largest/smallest value.

The ranking of chart elements

To find *unusual* values.

The distribution, ranking or magnitude of chart elements

# Why is someone reading your chart?

You have a story you want to tell

There's lots we can do to help guide the reader to understand your chart and follow the story you're telling. We'll cover some examples during this course.

---

The reader wants to see the data

Charts (and tables) are the best way to see the “big picture” of a dataset - a single value (eg mean) is kind of useless. Interactivity is really useful to allow readers to properly explore the dataset.

---

The reader has a preconception about the data

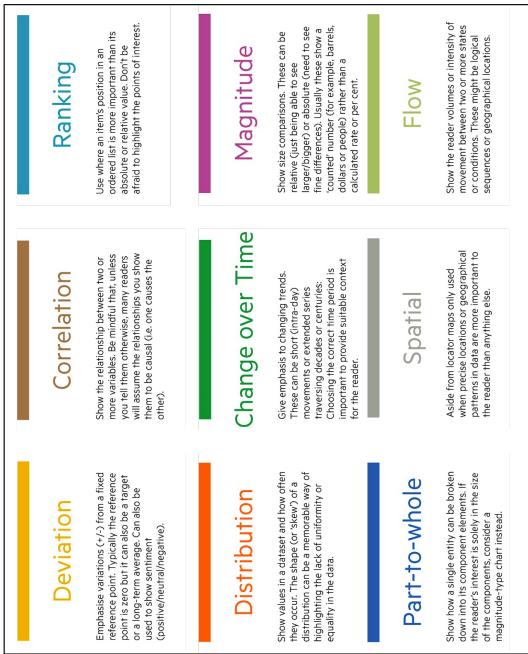
Readers might be approaching a chart biased with a particular theory about the data. We can do our best to make our charts easy to read and avoid common pitfalls.

# How do we choose a chart?

# Use data columns to choose charts

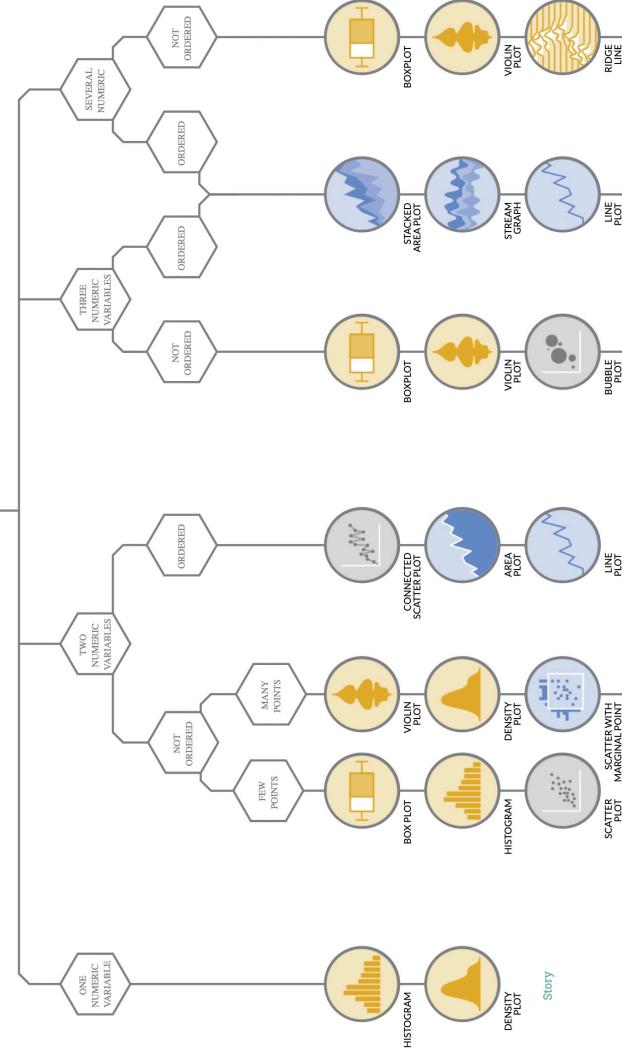
# Use your story to choose charts

```
> msleep
# A tibble: 83 × 11
  name   genus vore order cons...¹ sleep...² sleep...³ sleep...⁴ awake   brainwt
  <chr>  <chr> <chr> <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl>
  1 Cheetah Actin... carni Carni... lc    12.1  NA     11.9  NA      50
  2 Owl mo... Aotus omni Pr-in... NA    17     1.8  NA     7     0.0155  0.48
  3 Mounta... Aplo... herbi Rode... nt   14.4   2.4  NA     9.6  NA      1.35
  4 Greate... Blar... omni Sori... lc   14.9   2.3  0.133  9.1  0.00029  0.019
  5 Cow      Bos    herbi Arti... domest... 4     0.7   0.667  20    0.423   600
  6 Three-... Brad... herbi Pil... NA   14.4   2.2  0.767  9.6  NA      3.85
  7 Northe... Call... carni Carni... vu   8.7    1.4  0.383  15.3  NA      20.5
  8 Vesper ... Cat... NA      Rode... NA   7     NA     17   NA     0.045
  9 Dog      Canis carni Carni... domest... 10.1   2.9  0.333  13.9  0.07   14
  10 Roe de... Capr... herbi Arti... lc    3     NA     21   0.0982  14.8
  # ... with 73 more rows, and abbreviated variable names `¹conservation, `²sleep_total,
  #   `³sleep_rem, `⁴sleep_cycle
  # i Use `print(n = ...)` to see more rows
```



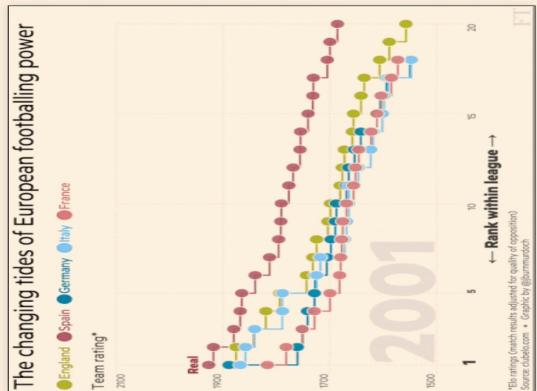
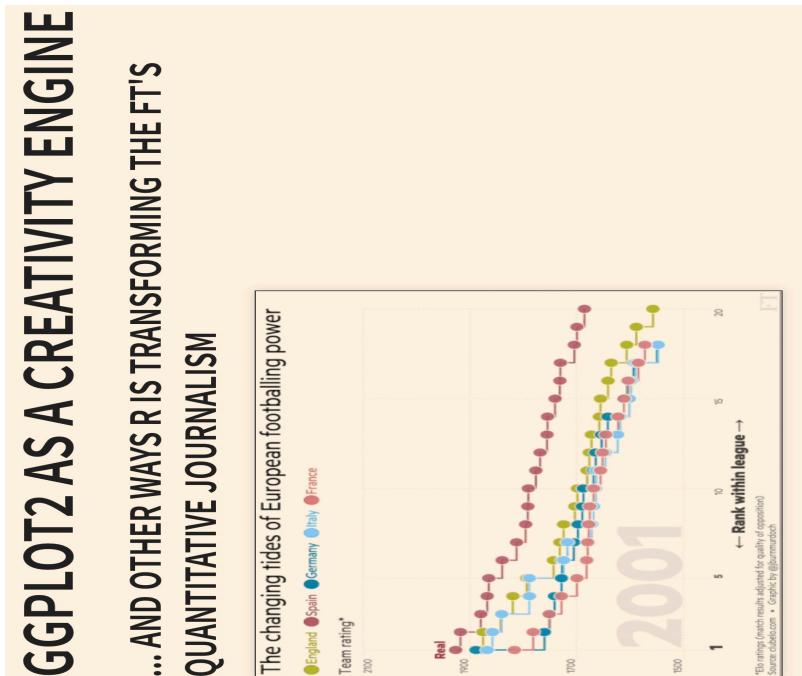
# data-to-viz.com

Numeric Categorical Num & Cat Maps Network Time series

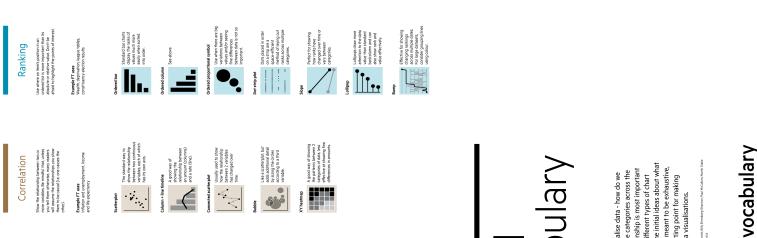
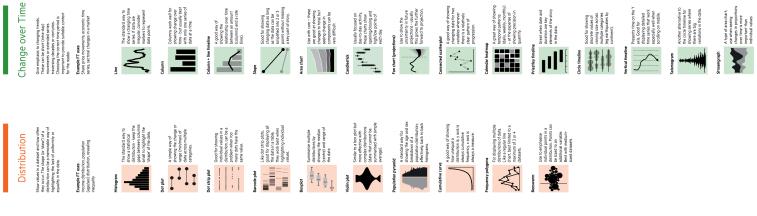
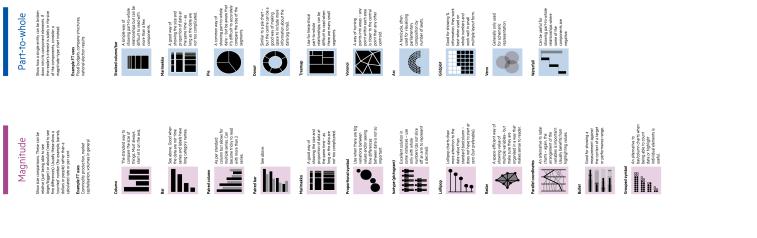
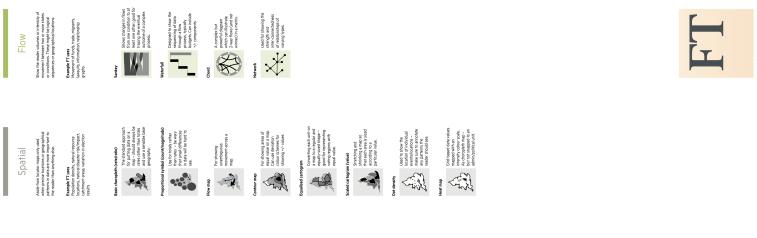


This site also provides simple to follow instructions for using {ggplot2} to build every single chart type you can find on the website.

[ft-interactive.github.io/vocabulary](https://ft-interactive.github.io/vocabulary)



International Law



## Visual vocabulary

There are many ways to visualize data - how do we know which one to pick? For the categories across the x-axis, it's likely that data relationships is most important to decide your chart type. When data is at the different types of chart within the category, to form some insights about what you're looking for. This is not meant to be exhaustive, nor is it a hard rule, but a useful starting point for creating informative and meaningful data visualizations.

[ft.com/vocabulary](http://ft.com/vocabulary)

# {ggplot2} for charts



# Task: Setup a new project

## SLIDE 1 OF 3

1. Create a new project called something like `week-4_dataviz.Rproj`
2. Add a new RMarkdown document called `ggplot2-notes.Rmd`

---

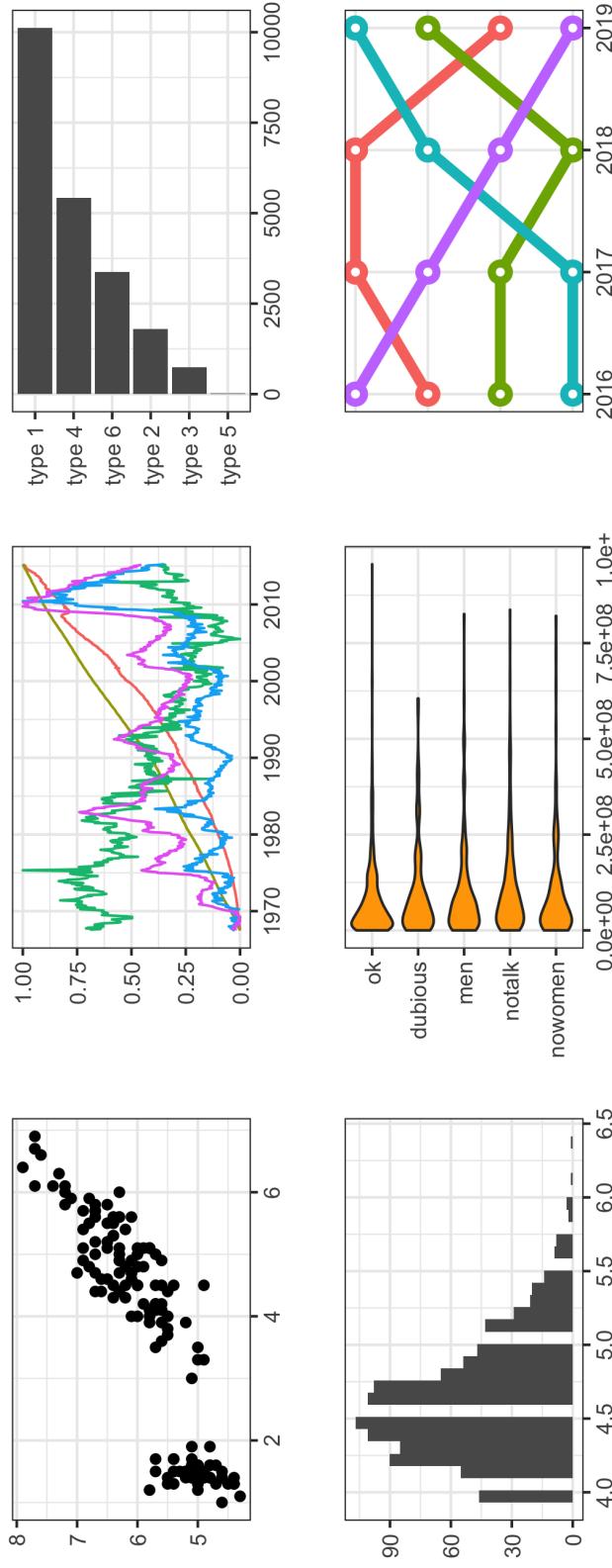
We're going to do some structured and unstructured code during today. During the workshop I'll be asking to you to create your own charts.

---

# ggplot2: A Grammar of Graphics

{ggplot2} is an incredibly powerful and flexible tool for building static dataviz.

We can build (almost)<sup>1</sup> any static chart we can conceive of.



[1] - Dual y-axis charts must be transformations of one another (**for good reasons**)

# Building blocks of a `{ggplot2}` chart



Aesthetics



Geoms



Scales



Guides



Theme

## ■ Aesthetics

Aesthetics are used to create **mappings** between columns in our datasets and the coordinate systems of our chart:

```
1 msleep %>%
2   ggplot() +
3   aes(
4     x = sleep_total,
5     y = sleep_rem,
6     colour = vore
7   )
```

{ggplot2} uses *tidy evaluation* to allow us to use bare column names in our code.

## ■ Aesthetics

### Where is aes() placed? What it does

---

Inside <code>ggplot()</code> or on its own	Sets the aesthetics for the entire <code>{ggplot2}</code> object.
--	---

These could be considered the *coordinate system aes()*

---

Inside <code>geom_*</code> ()	Sets aesthetics for a specific geom <i>within</i> the existing coordinate system aes() for the <code>{ggplot2}</code> object.
-------------------------------	---

These should be considered *geom specific aes()*

Geoms use the aesthetics to add layers to our charts.

```
1 msleep %>%
2   ggplot() +
3   aes(
4     x = sleep_total,
5     y = sleep_rem,
6     colour = voro
7   ) +
8   geom_point()
```

There are 50+ geoms baked into the `{ggplot2}` package.

## ▲ Geoms

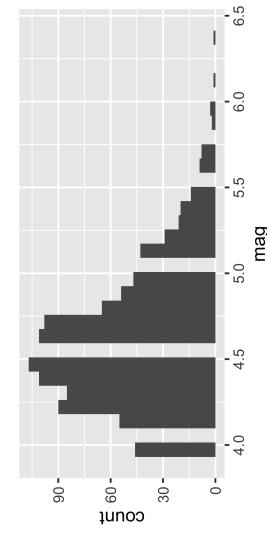
```
geom_abline(), geom_area(), geom_bar(), geom_bind2d(), geom_blank(), geom_boxplot(),
geom_col(), geom_contour(), geom_countour_filled(), geom_count(), geom_crossbar(),
geom_curve(), geom_density(), geom_density2d(), geom_density2d_filled(),
geom_density2d(), geom_dotplot(), geom_errorbar(),
geom_errorbarh(), geom_freqpoly(), geom_function(), geom_hex(), geom_histogram(),
geom_hline(), geom_jitter(), geom_label(), geom_line(), geom_linerange(), geom_map(),
geom_path(), geom_point(), geom_pointrange(), geom_polygon(), geom_qq(),
geom_qq_line(), geom_quantile(), geom_raster(), geom_rect(), geom_ribbon(),
geom_rug(), geom_segment(), geom_sf(), geom_sf_label(), geom_sf_text(),
geom_smooth(), geom_spoke(), geom_step(), geom_text(), geom_tile(), geom_violin(),
geom_vline()
```

As we'll see later, there are many `{ggplot2}` extension packages that add even more geoms to the mix.

# Some geoms are built from others (I)

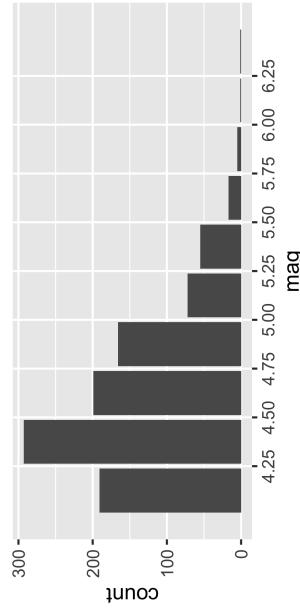
`geom_histogram()` has clever tricks to make useful histograms

```
1  ggplot(quakes, aes(mag)) +  
2  geom_histogram()  
3  scale_x_binned()
```



It's built by calling `geom_bar()`

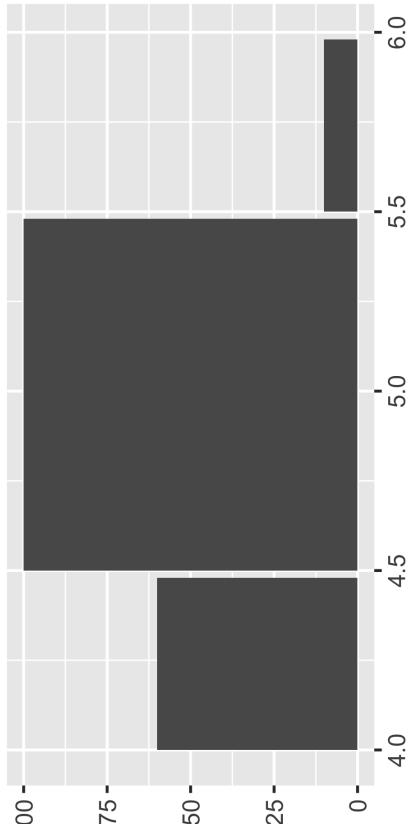
```
1  ggplot(quakes, aes(mag)) +  
2  geom_bar() +  
3  scale_x_binned()
```



# Some geoms are built from others (II)

But `geom_bar()` itself is built from `geom_rect()`.

```
1 rect_data <- tribble(
2   ~x_min, ~x_max, ~y_min, ~y_max,
3   4, 4.48, 0, 60,
4   4.5, 5.48, 0, 100,
5   5.5, 5.98, 0, 10
6 )
7 rect_data %>%
8   ggplot() +
9     geom_rect(aes(xmin = x_min,
10                  xmax = x_max,
11                  ymin = y_min,
12                  ymax = y_max)) +
13   theme_gray(base_size = 25)
```



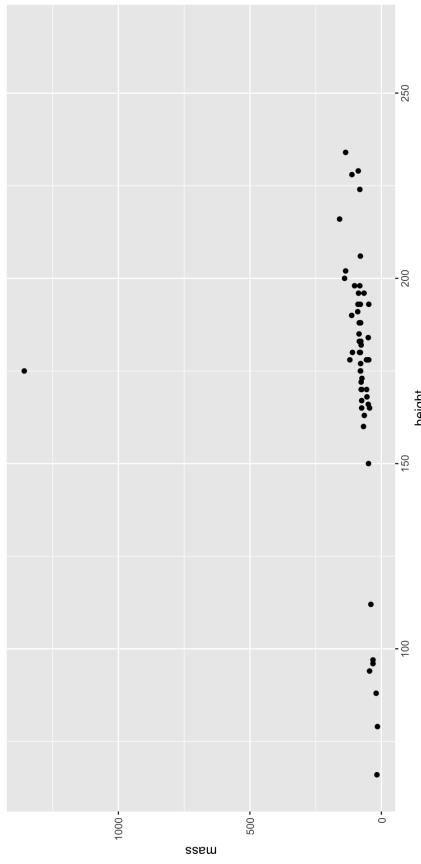
| There are **8 primitives** from which all other geoms are built:

`geom_blank()`, `geom_point()`, `geom_polygon()`, `geom_rect()`,  
`geom_ribbon()`, `geom_segment()`, `geom_text()`

# All geoms have x and y aesthetics

These tell the geom where it needs to be drawn:

```
1 starwars %>%
2   ggplot() +
3   aes(x = height,
4        y = mass) +
5   geom_point()
```



# Some geoms need more than just x and y

Let's `geom_segment()` to visualise some of the eras of the dinosaurs:

```
1 dinosaurs <- tribble(  
2   ~period, ~start, ~end,  
3   "Triassic Period", -251e6, -225e6,  
4   "Late Triassic Period", -225e6, -200e6,  
5   "Jurassic Period", -200e6, -150e6,  
6   "Late Jurassic Period", -150e6, -145e6  
7 )
```

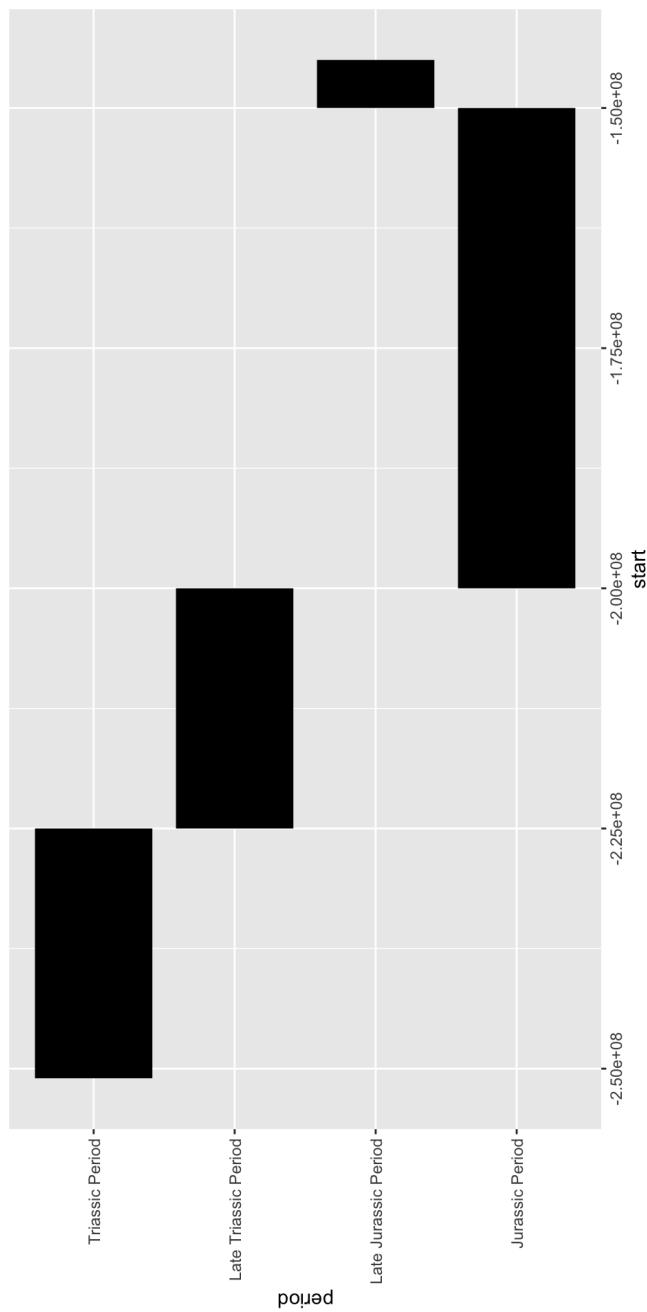
To build this chart we need to specify **all** of the following: x, xend, y and yend.

# Use `size` to affect geom size

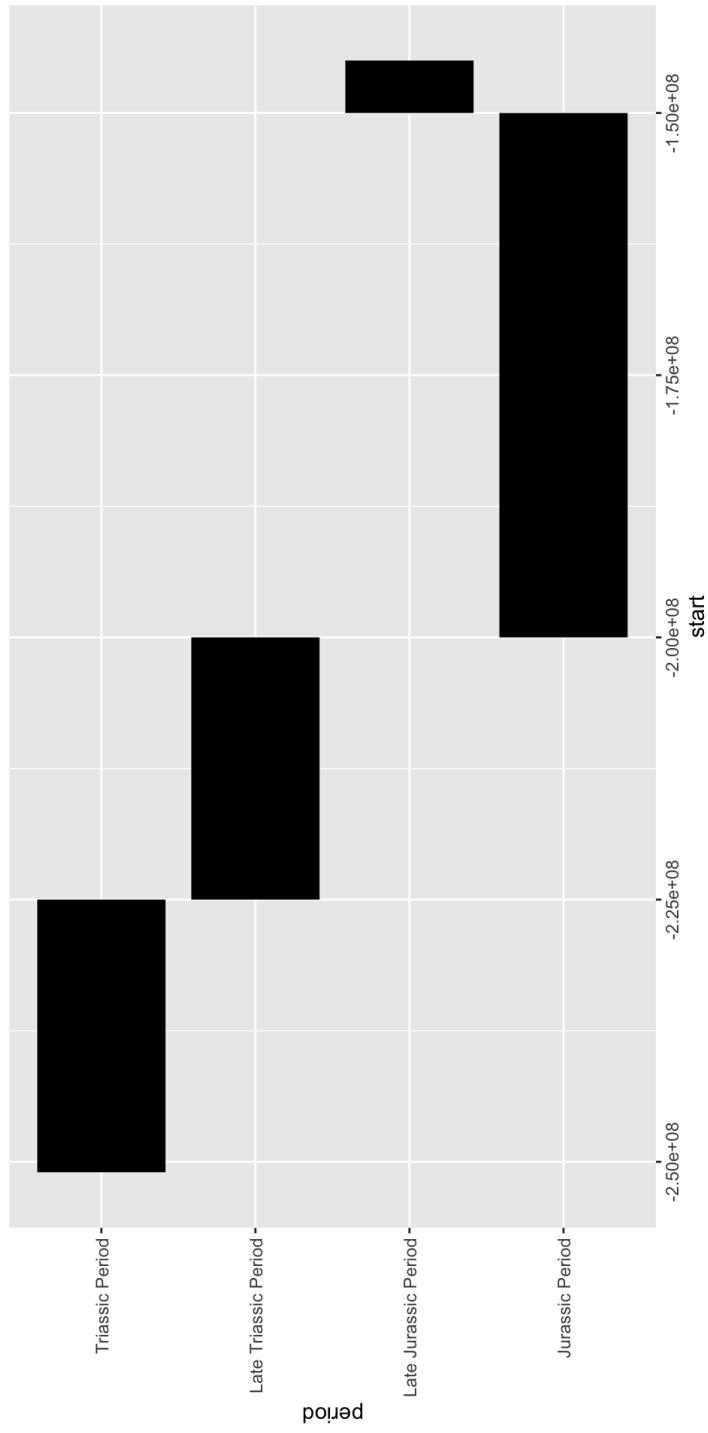
In many charts we want geoms to be **thicker**, **bigger** or *just be more prominent.*

Timeline (or Gantt charts) are good examples of this. We want the segments to be thicker to improve the readability of the chart - this comes down to the `size` aesthetic.

```
1 dinosaurs %>%  
2 ggplot() +  
3 aes(x = start, xend = end,  
4     y = period, yend = period) +  
5 geom_segment(size = 30)
```



# Out of order dinosaurs



This is still a **bad chart**.

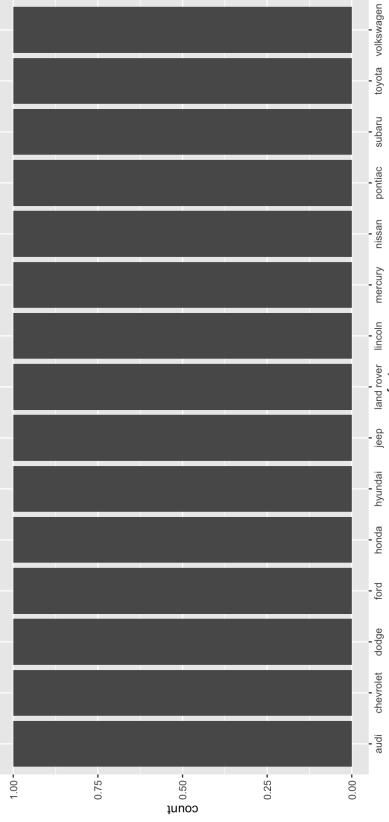
The eras are not ordered in geological time, instead they're ordered (reverse) alphabetically.

To control the order of things in `{ggplot2}` charts we must use **factors** - which are picked up by the **scales**.

# Some geoms are designed to save time

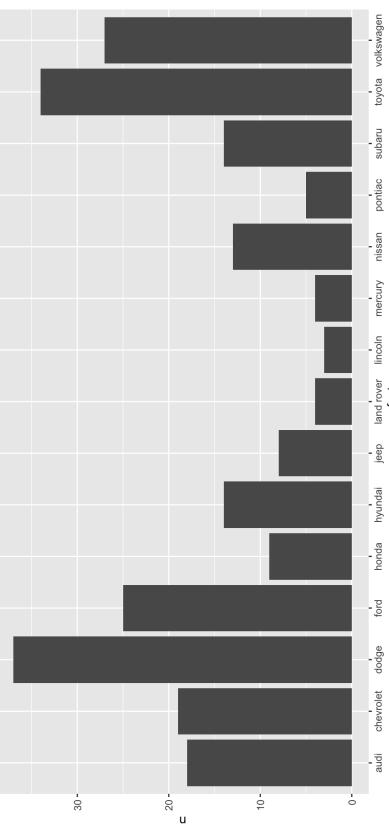
`geom_bar()` defaults to counting instances of a variable.

```
1 mpg %>%
2   count(manufacturer) %>%
3   ggplot() +
4   geom_bar(aes(manufacturer))
```



`geom_col()` uses a column to dictate the length of bars.

```
1 mpg %>%
2   count(manufacturer) %>%
3   ggplot() +
4   geom_col(aes(x = manufacturer, y = n))
```

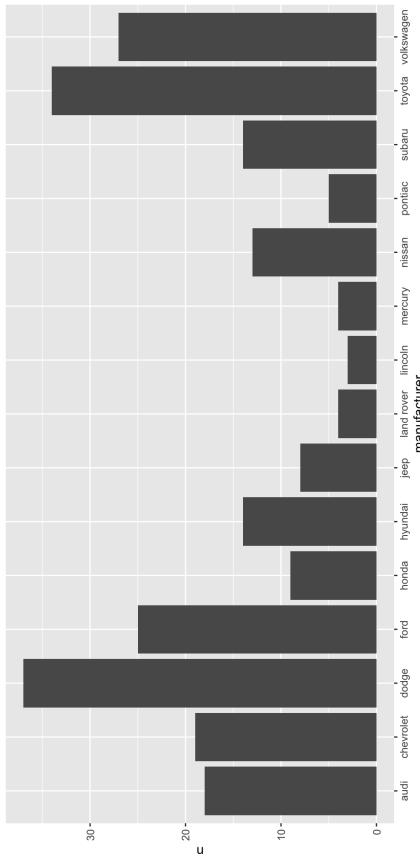


# Some geoms depend on stat functions

The `geom_bar()` function has a `stat` argument with the default value of "count".

We can force the geom to behave like `geom_col()` by changing the stat:

```
1 mpg %>%
2   count(manufacturer) %>%
3   ggplot() +
4   geom_bar(aes(x = manufacturer,
5               y = n),
6             stat = "identity")
```



All of the goodness from the `stat` argument comes from the `stat_identity()` and `stat_count()` functions.

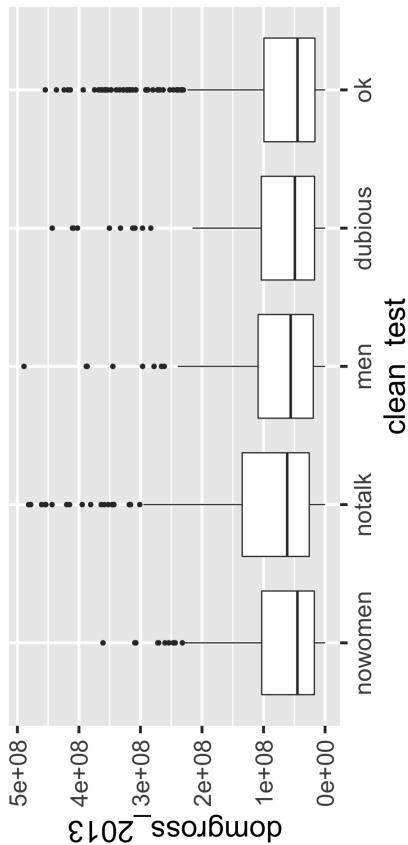
If you're building a complex chart it might be useful to directly call a `stat_()` function.

# Position things to resolve overlapping

## (1)

Box and whisker diagrams hide a lot of detail

```
1 bechdel %>%
2   filter(complete.cases(.),
3         domgross_2013 < 0.5e9) %>%
4   ggplot(aes(clean_test,
5             domgross_2013)) +
6   geom_boxplot() +
7   theme_gray(base_size = 25)
```



Let's add the data points to this chart with `geom_point()` and look at the position argument.

# Position things to resolve overlapping

## (II)

The position argument can also be used to create three different types of bar chart:

- “stack” creates a stacked bar chart
- “fill” creates a proportional bar chart
- “dodge” creates a grouped bar chart

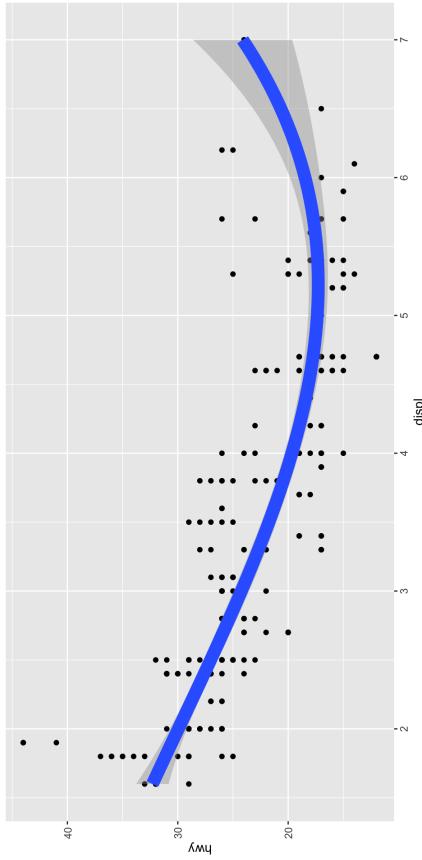
Let's create all 3 of these for the following dataset:

```
1 gss_cat %>%
2   count(relig, marital)

# A tibble: 78 × 3
  relig    marital      n
  <fct>   <fct>     <int>
  1 No answer No answer     4
  2 No answer Never married 22
  3 No answer Separated      3
  4 No answer Divorced     13
  5 No answer Widowed       7
  6 No answer Married      44
  7 Don't know Never married 6
  8 Don't know Separated      3
  9 Don't know Divorced      1
 10 Don't know Married        5
# ... with 68 more rows
```

# Geom layers are placed on top of one another

```
1 ggplot(mpg,
2       aes(displ, hwy)) +
3     geom_point() +
4     geom_smooth(method = lm,
5                 formula = y ~ splines::bs(x, 3),
6                 size = 5)
```



The `geom_smooth()` line is hiding data points.

We could either swap the order of these geoms or change the `alpha` aesthetic.

Scales determine the appearance of an aesthetic within the chart, including:

- Axes labels and breaks
- Colours used for `colour` and `fill` aesthetics

```
1 mslleep %>%
2   ggplot() +
3   aes(
4     x = sleep_total,
5     y = sleep_rem,
6     colour = vore
7   ) +
8   geom_point() +
9   scale_colour_manual(
10    values = c("carni" = "#c03728",
11      "omni" = "#fd8f24",
12      "insecti" = "#f5c04a",
13      "herbi" = "#919c4c",
14      "NA" = "#e68c7c")
15  )
```

## Δ Scales

Scales also determine the **order** in which elements are shown in a chart.

To change the order of discrete/categorical columns we need to use **factors**.

# Scales and {scales}

{ggplot2} uses the {scales} package under the hood to build all of the scales that we see - including continuous and discrete scales.

The {scales} package also contains many utility functions that are useful for us to format our axes and other scales.

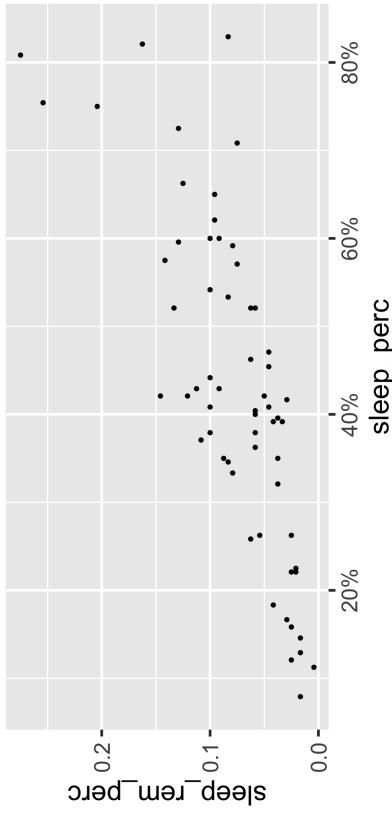
We can either load the {scales} package itself or call functions specifically with  
`scales::label_percent()`

# {scales} and depreciation (I)

Until recently the way we'd use `{scales}` would be as follows

```
1 percent(c(0.3, 0.5, 0.6))  
[1] "30%" "50%" "60%"
```

```
1 msleep %>%  
2   mutate(sleep_perc = sleep_total / 24,  
3         sleep_rem_perc = sleep_rem / 24) %>%  
4   ggplot() +  
5   aes(x = sleep_perc,  
6        y = sleep_rem_perc) +  
7   geom_point() +  
8   scale_x_continuous(label = percent_format()) +  
9   theme_gray(base_size = 24)
```



There was a function called `percent(x)` for formatting a vector of values `x` and `percent_format()` for modifying the appearance of percentages in a `{ggplot2}` chart.

# {scales} and depreciation (II)

These functions have now been deprecated. This means there are new alternatives to these functions.

Deprecation is a fact of life in software development. But the details of how things are deprecated are variable.

Sometimes things are deprecated with the intention of removing them in the future. Other times, the deprecated functions will continue to exist far into the future.

It seems like the intention is for these functions to continue to work into the future. But they might be removed in several years time.

---

Let's use the new approach for formatting scales so that you can read modern documentation and so you're not learning deprecated functions.

# {scales} and depreciation (!!!)

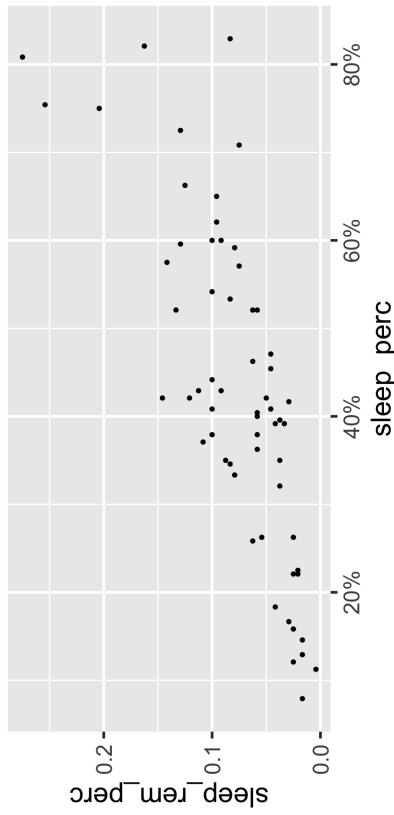
We now use `label_percent()` for both types of operation.

```
1  label_percent()(c(0.3, 0.5, 0.6))  
[1] "30%" "50%" "60%"
```

This is known as a function factory.

Function factories are cool. But I wish you didn't have to learn this syntax.

```
1  msleep %>%  
2  mutate(sleep_perc = sleep_total / 24,  
3        sleep_rem_perc = sleep_rem / 24) %>%  
4  ggplot() +  
5  aes(x = sleep_perc,  
6       y = sleep_rem_perc) +  
7  geom_point() +  
8  scale_x_continuous(label = label_percent()) +  
9  theme_gray(base_size = 24)
```

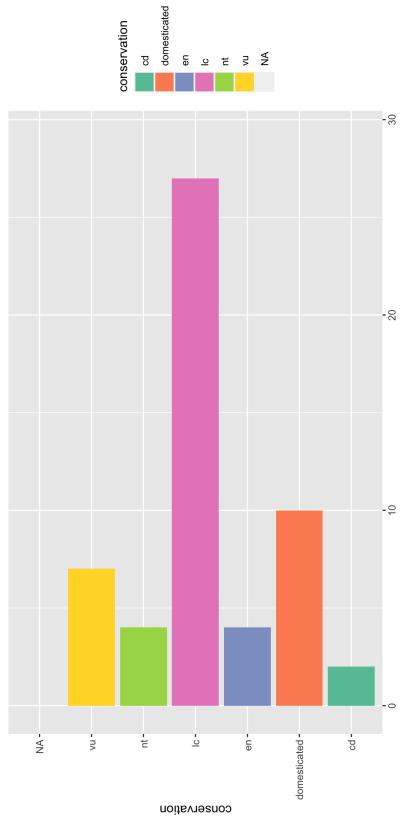


# {scales} and colours ()

There are many built-in colour palettes in {scales} - let me introduce two families of palettes.

# {scales} and colours (II)

There are many built-in colour palettes in {scales} - let me introduce two families of palettes.



There are some pretty good palettes for discrete/categorical variables in this family of palettes.

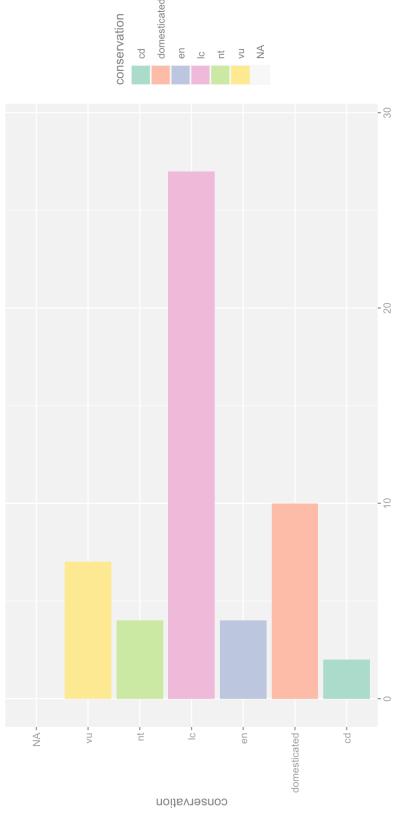
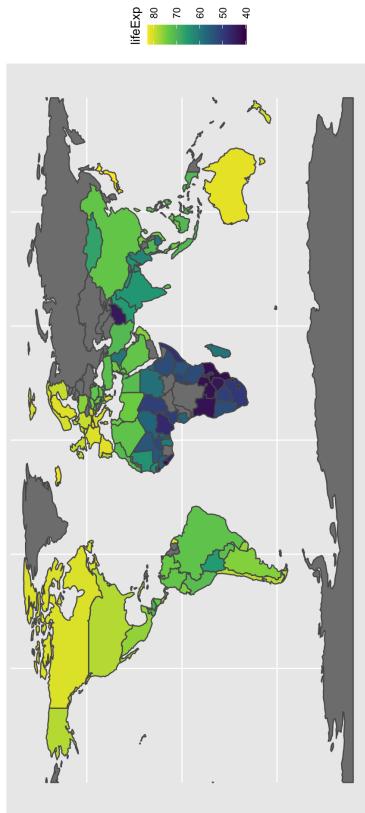
# {scales} and colours (III)

There are many built-in colour palettes in {scales} - let me introduce two families of palettes.

But for continuous variables I strongly recommend using the **viridis** family of palettes.

These are designed to be both perceptually uniform and to work for folks with colour blindness.

```
1 countries110 %>%
2     st_as_sf() %>%
3     left_join(filter(gapminder, year == 2007),
4                by = c("name" = "country")) %>%
5     ggplot() +
6     geom_sf(aes(fill = lifeExp)) +
7     scale_fill_viridis_c()
```

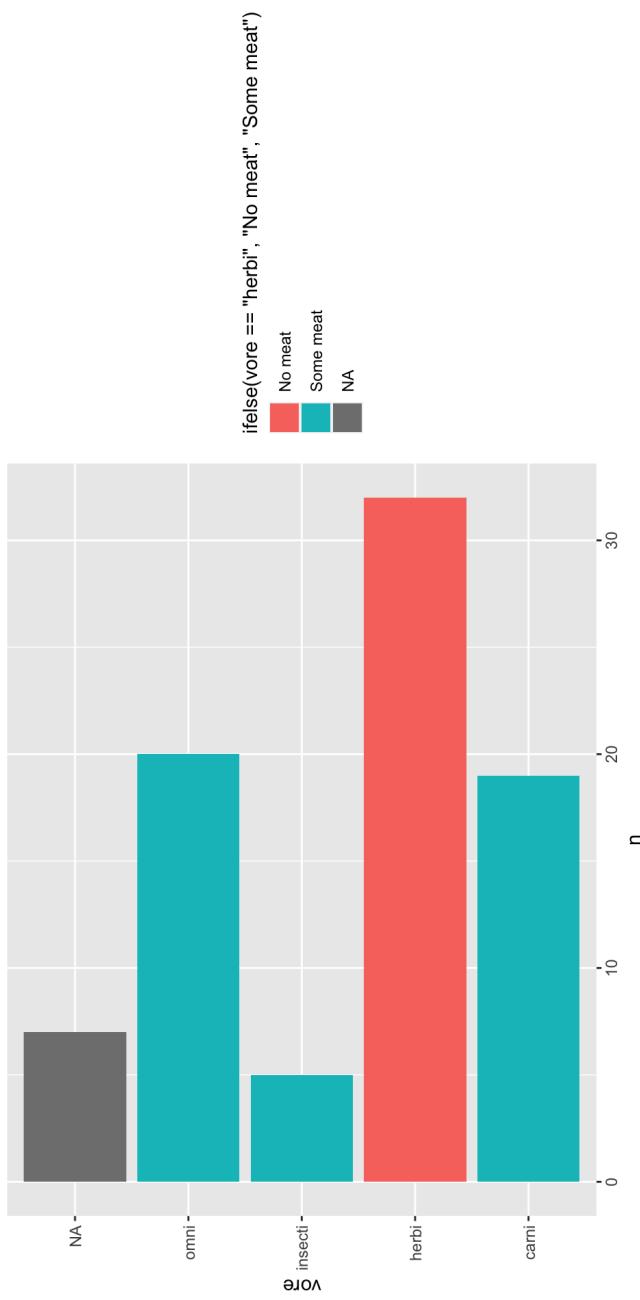


There are some pretty good palettes for discrete/categorical variables in this family of palettes.

# Setting custom colours ()

One of the first frustrations people find with {ggplot2} is setting our own custom colours, eg in this chart:

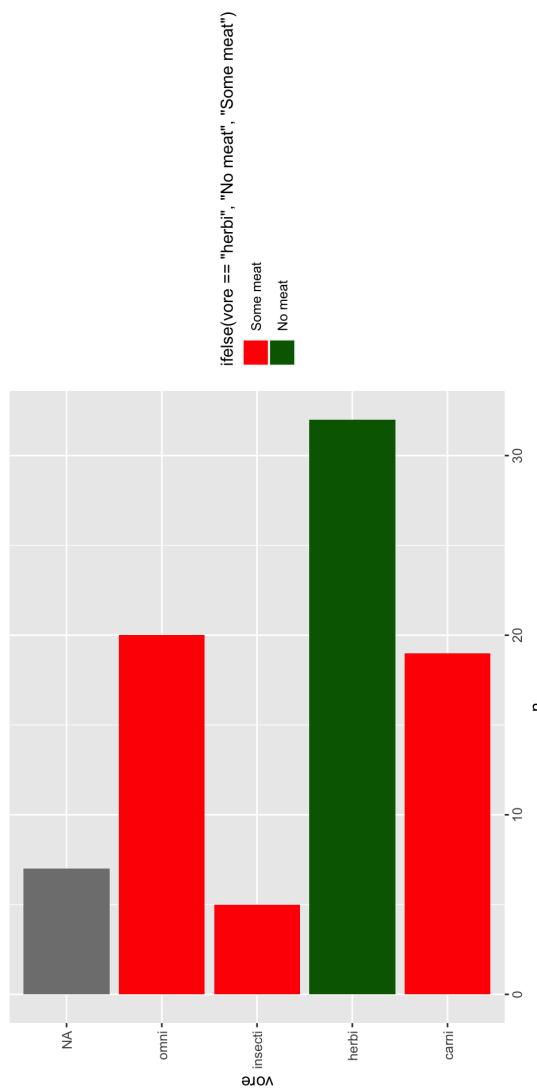
```
1 msleep %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
7   geom_col()
```



# Setting custom colours (II)

We need to use `scale_fill_manual()`

```
1 mslleep %>%
  2   count(vore) %>%
  3     ggplot() +
  4       aes(x = n,
  5         y = vore,
  6         fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
  7         geom_col() +
  8         scale_fill_manual(values = c("Some meat" = "red",
  9           "No meat" = "darkgreen"))
```



We'll come back to this chart in the section on guides().

# {**scales**} and **factors** (I)

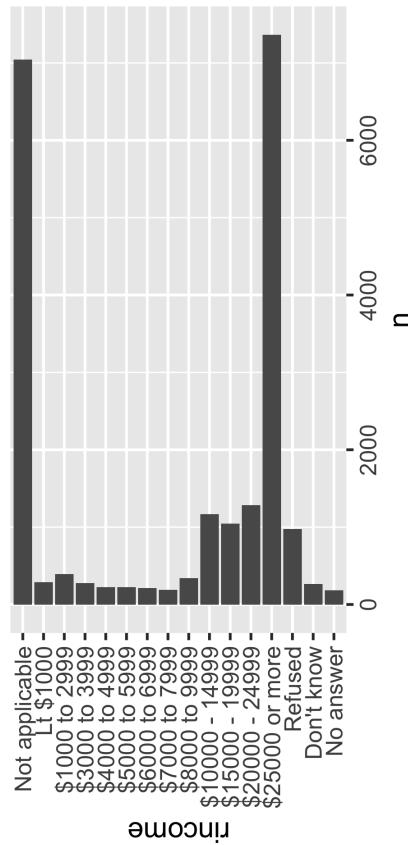
Factors are R's categorical data type. They allow us to create a variable with fixed values (**levels**) and to set the *order* of those levels.

Let's look at a pre-existing dataset with factors:

```
1 gss_cat %>%
  2 head() %>%
  3 pull(rincome)
```

```
[1] $8000 to 9999 $8000 to 9999 Not applicable Not applicable Not applicable
[6] $20000 - 24999
16 Levels: No answer Don't know Refused $25000 or more ... Not applicable
```

```
1 gss_cat %>%
  2 count(rincome) %>%
  3 ggplot() +
  4 aes(x = n,
  5   y = rincome) +
  6 geom_col() +
  7 theme_gray(base_size = 24)
```



# {**scales**} and **factors** (II)

The base R tools for creating and manipulating factors are messy and frustrating to use.

We're going to use the {forcats} package which is loaded when we run  
`library(tidyverse)`.

Almost all of the functions begin with `fct_*` to let you know we're dealing with factors.

# factors and msleep

Let's think of the different ways we could order this dataset:

```
1 msleep %>%
  2 count(vore)
```

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

## Count order

In this ordering we will arrange the `vore` column according to values in the `n` column.

## Canonical order

In this ordering we'll arrange the `vore` column from the diet with the most meat to the least meat.

This is usually what we want in count bar charts.

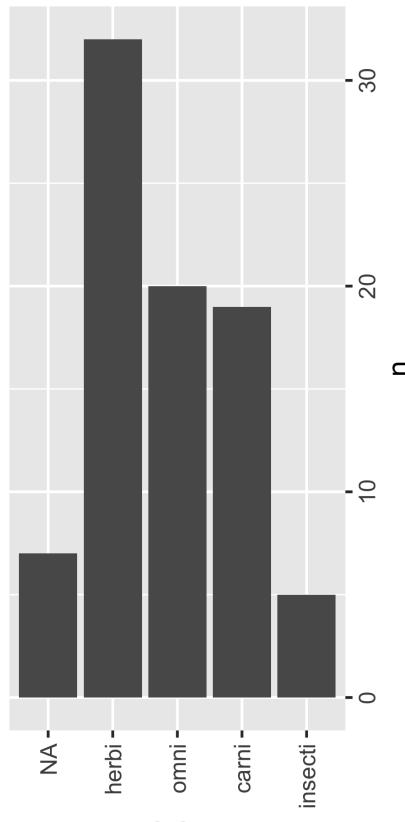
This is usually what we want in visualising survey datasets,

- eg Strong disagree, Disagree, Neither agree or disagree, Agree, Strong Agree

# msleep factor: count order (I)

We use `fct_reorder()` to order a factor by another column.

```
1 msleep %>%
2   count(vore) %>%
3   mutate(vore = fct_reorder(vore, n)) %>%
4   ggplot() +
5   aes(x = n,
6        y = vore) +
7   geom_col() +
8   theme_gray(base_size = 24)
```

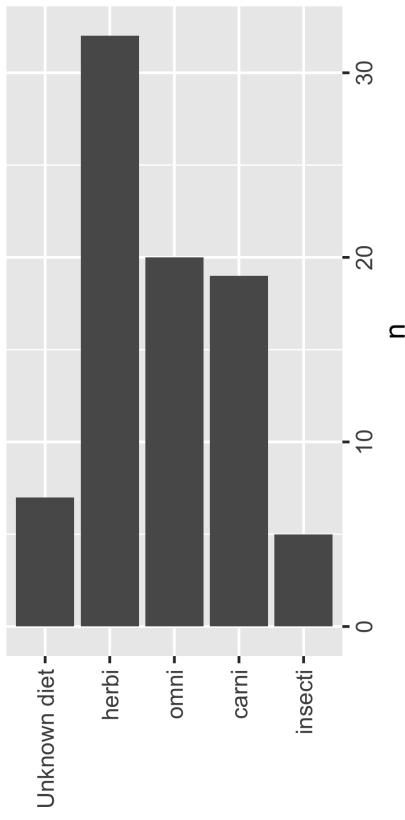


But what about the NA values? What should we do?

# msleep factor: count order (III)

We can replace NA values nicely with `fct_explicit_na()`

```
1 msleep %>%
2   count(vore) %>%
3   mutate(vore = fct_reorder(vore, n),
4         vore = fct_explicit_na(vore, "Unknown diet"))
5 ggplot() +
6   aes(x = n,
7       y = vore) +
8   geom_col() +
9   theme_gray(base_size = 24)
```

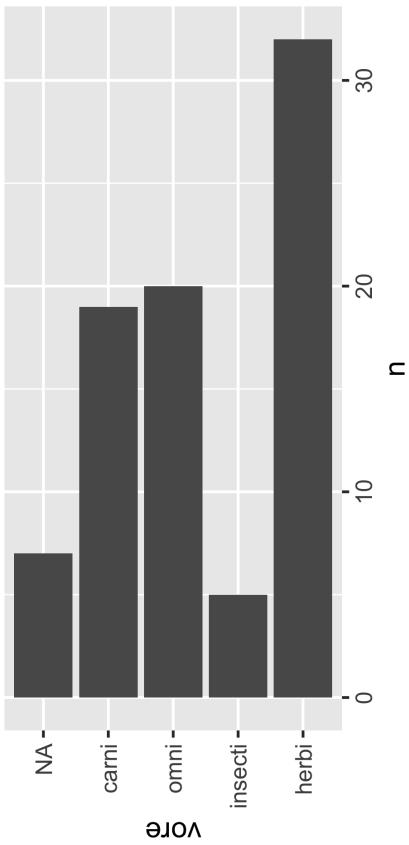


Let's come back to moving the position of the NA level.

# msleep factor: canonical order ()

To set our own canonical order we use `fct_relevel()` and provide a vector with our preferred order.

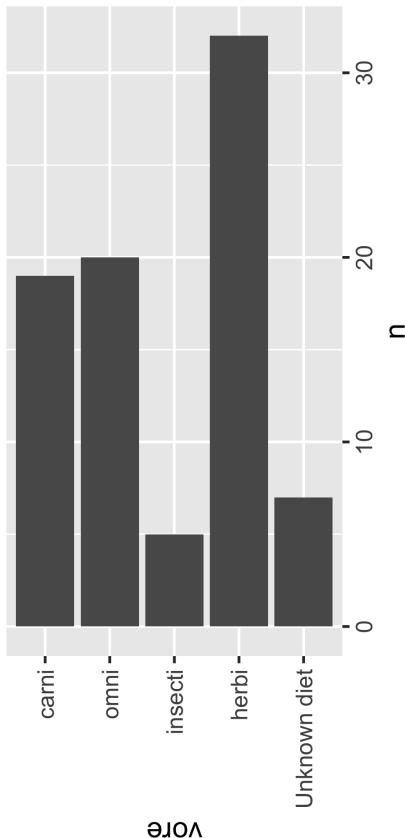
```
1 order_vore <- c("carni", "omni", "insecti", "herbi")
2
3 msleep %>%
4   count(vore) %>%
5   mutate(vore = fct_relevel(vore, order_vore),
6         vore = fct_rev(vore)) %>%
7   ggplot() +
8   aes(x = n,
9        y = vore) +
10  geom_col() +
11  theme_gray(base_size = 24)
```



# msleep factor: canonical order (II)

We can also use `fct_relevel()` to modify the position of a specific

```
1 msleep %>%
2   count(vore) %>%
3   mutate(vore = fct_relevel(vore, order_vore),
4         vore = fct_rev(vore),
5         vore = fct_explicit_na(vore, "Unknown diet"),
6         vore = fct_relevel(vore, "Unknown diet", after = 1))
7 ggplot() +
8   aes(x = n,
9        y = vore) +
10  geom_col() +
11  theme_gray(base_size = 24)
```





# Task: Global Burden of Disease factors

## SLIDE 1 OF 3

These are the same steps you've repeated before

1. Add a sub-folder to your project called **data**
2. Inside of the **data** folder add a script called **obtain-data.R**
3. Add this code
5. Run the 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")
```



# Task: Global Burden of Disease factors

## SLIDE 2 OF 3

1. Add a new heading for the GBD Dataset to your .Rmd

2. Filter the dataset as follows:

- Most recent year
- `location_name` starts with “World Bank”
- `metric_name` is “Number”
- `cause_name` is “Injuries”

3. Select only these columns

- `location_name, cause_name, val`

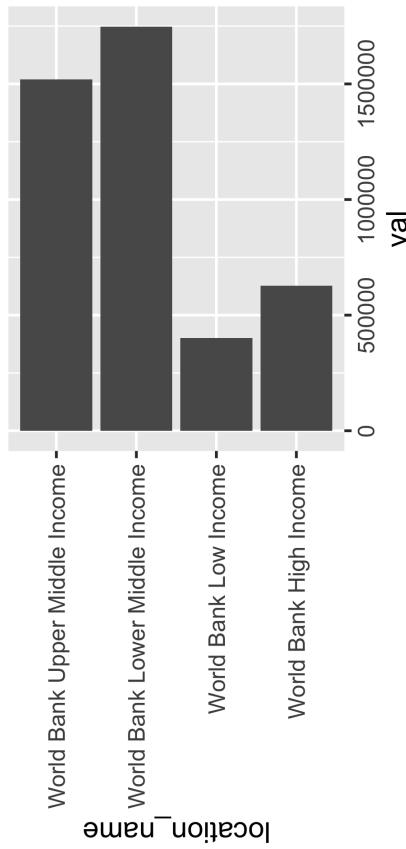


# Task: Global Burden of Disease factors

## SLIDE 2 OF 3

Create two versions of this chart:

- Bars are ordered by their size
- Bars are ordered from “World Bank High Income” to “World Bank Low Income”



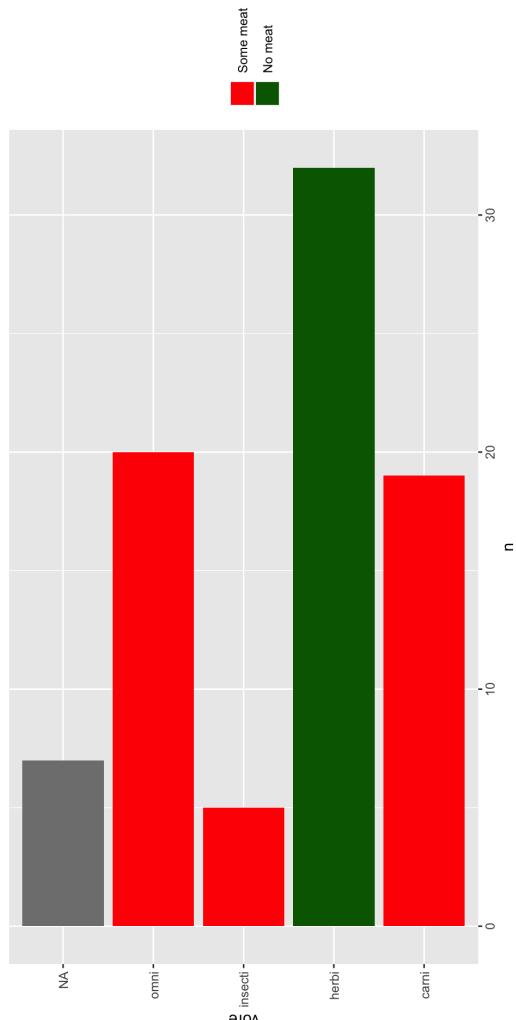
```
1 gdb_injuries %>%
2   ggplot() +
3   aes(x = val,
4        y = location_name) +
5   geom_col() +
6   theme_gray(base_size = 24)
```

We (kind of) use “guides” and “legends” interchangeably in `{ggplot2}`.

```
1 ggplot() +
2   geom_line(show.legend = FALSE) +
3   guides(alpha = guide_legend() +
4     theme(legend.position = "bottom"))
```

Guides can **only** be created through a corresponding aesthetic and scale.

## » Guides



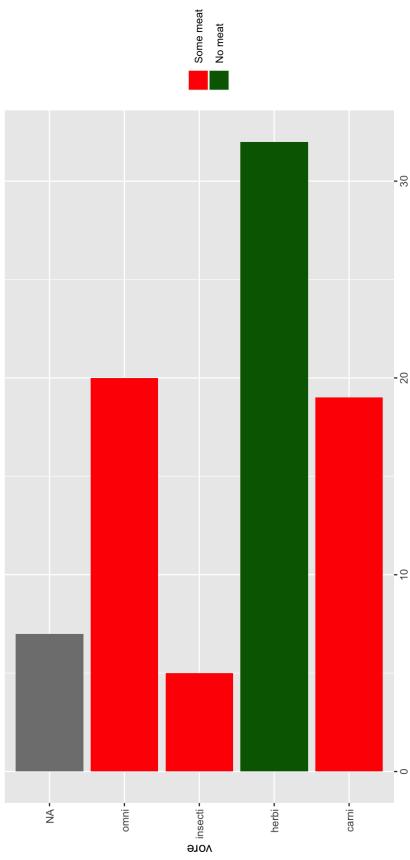
# “Manual legends” (I)

Sometimes we want to add additional legend items - usually for NA values, and particularly for maps.

Let's continue with this chart from before:

```
1 msleep %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi",
7                   "No meat",
8                   "Some meat")) +
9   geom_col() +
10  scale_fill_manual(values = c("Some meat" = "red",
11                      "No meat" = "darkgreen",
12                      name = ""))

```

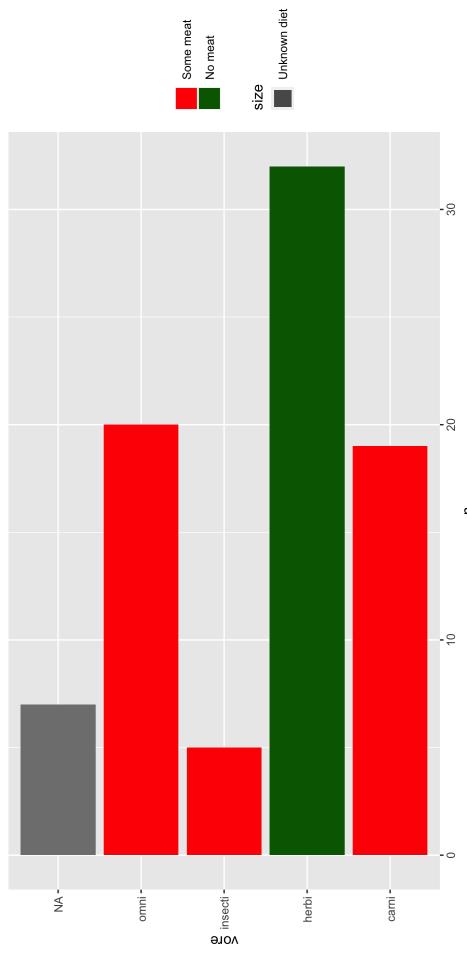


# “Manual legends” (III)

We need to choose an aesthetic that works for `geom_col()` but we’re not using elsewhere in the chart.

This will change depending on your chart. In this instance we can use `size`

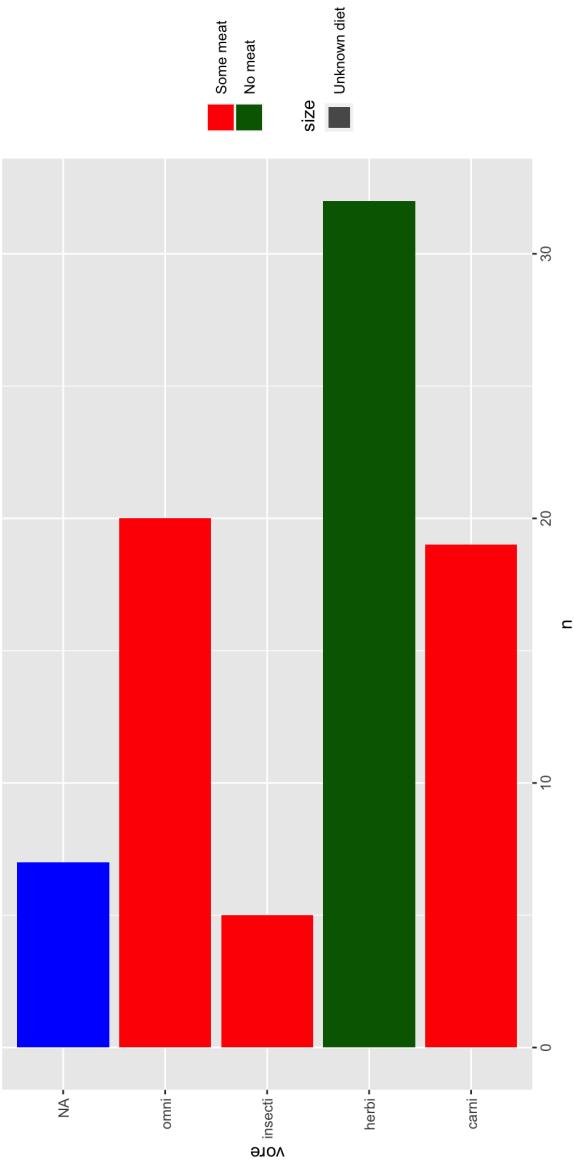
```
1 mslrep %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
7   geom_col(aes(size = "Unknown diet")) +
8   scale_fill_manual(values = c("Some meat" = "red",
9                             "No meat" = "darkgreen"),
10                      name = "")
```



# ‘Manual legends’ (III)

We now set the `na.value` colour for the original `scale_fill_manual()` scale

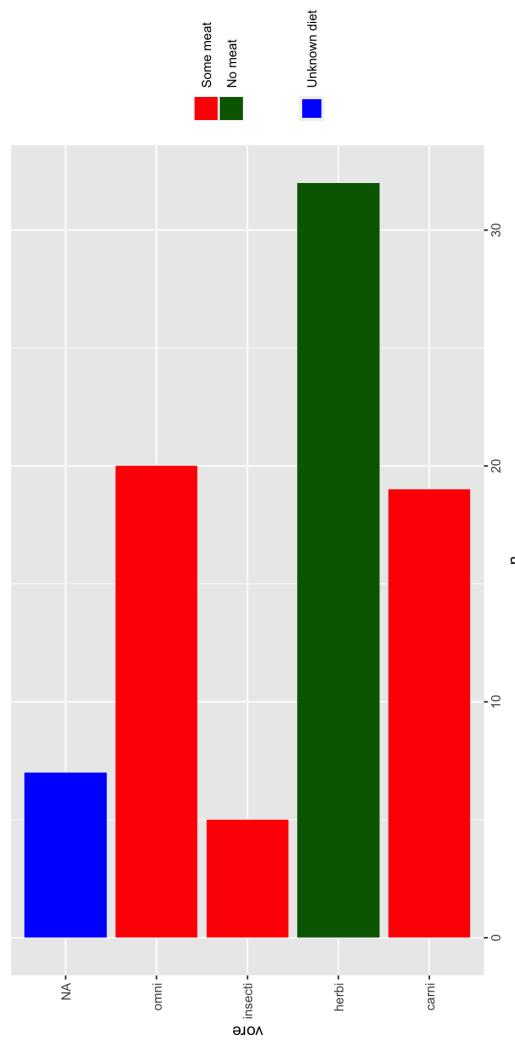
```
1 msleep %>%
 2   count(vore) %>%
 3   ggplot() +
 4   aes(x = n,
 5     Y = vore,
 6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
 7   geom_col(aes(size = "Unknown diet")) +
 8   scale_fill_manual(values = c("Some meat" = "red",
 9                             "No meat" = "darkgreen"),
10                     name = "",
11                     na.value = "blue")
```



# “Manual legends” (IV)

Next we use the `guides()` function to override the values for the `size` legend

```
1 msleep %>%
 2   count(vore) %>%
 3   ggplot() +
 4   aes(x = n,
 5     Y = vore,
 6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
 7   geom_col(aes(size = "Unknown diet")) +
 8   scale_fill_manual(values = c("Some meat" = "red",
 9     "No meat" = "darkgreen"),
10   name = "",
11   na.value = "blue") +
12   guides(size = guide_legend(title = "",
13   override.aes = list(fill = "blue")))
```



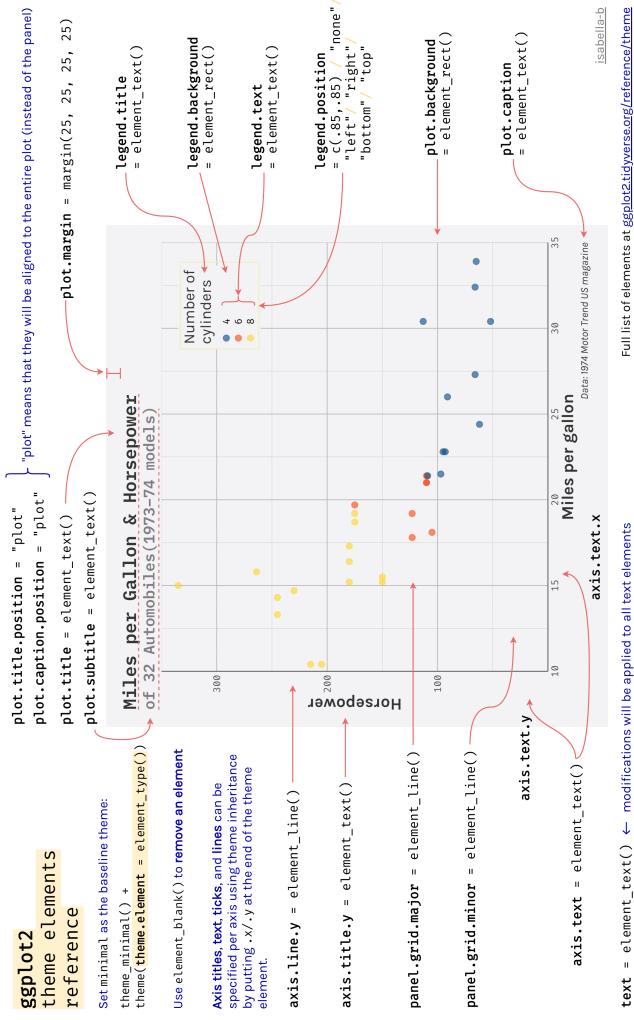
# More about guides()

If we want to modify the size of legend items we have two choices:

- `guides(fill = guide_colourbar(barwidth = 0.5, barheight = 10))`
- ... or to set the sizes in the `theme()`.

There are over 92 arguments to the `theme()` function for controlling chart appearance.

Remembering them all is challenging - I usually google them! Or use guides like this one:



Source: <https://bookdown.org/alapo/learnr/data-visualisation.html>

text = element\_text() ← modifications will be applied to all text elements  
Full list of elements at [https://plot2.ridvansorg/reference/theme\\_isabelina-b](https://plot2.ridvansorg/reference/theme_isabelina-b)

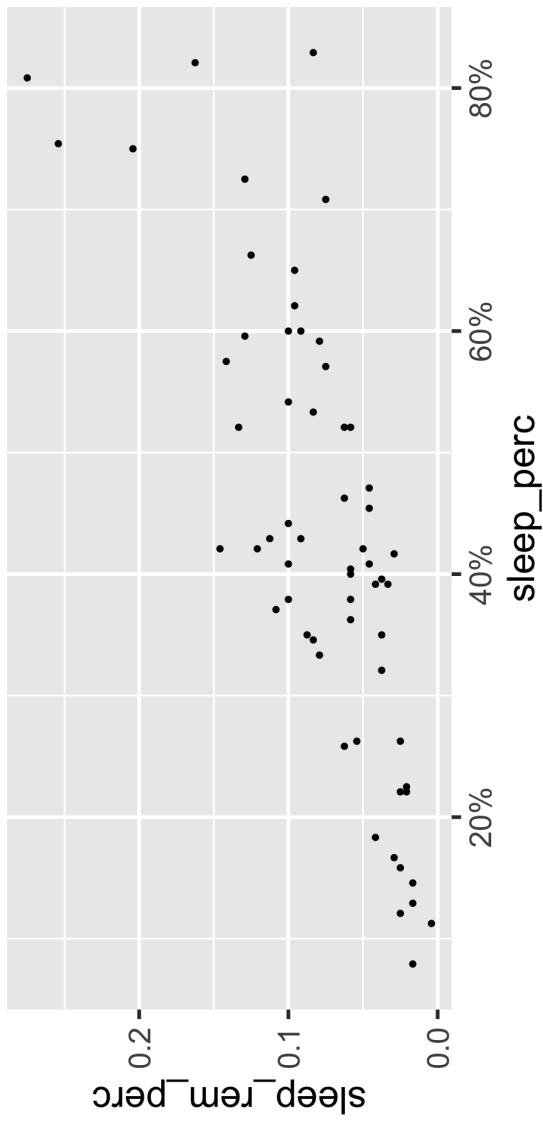


# What themes are there? (I)

{ggplot2} has several built-in themes. They have several arguments for quickly customising them.

I've been using the default `theme_gray()` to change text size in charts.

```
1 mslleep %>%
  2   mutate(sleep_perc = sleep_total / 24,
  3         sleep_rem_perc = sleep_rem / 24) %>%
  4   ggplot() +
  5   aes(x = sleep_perc,
  6        y = sleep_rem_perc) +
  7   geom_point() +
  8   scale_x_continuous(label = label_percent()) +
  9   theme_gray(base_size = 24)
```



# What themes are there? (II)

The {ggthemes} package contains lots of really useful - and beautiful - themes.

It's recommended that you choose a theme close to what to want and then customise it.

```
1 theme_fivethirtyeight() +  
2   theme(panel.grid.major = element_line(colour = "red"))
```

# **element\_\*()** functions in legends

Most of the legend arguments expect one of these functions:

- `element_line()`
- `element_text()`
- `element_rect()`

Or `element_black()` if you want to remove a theme element.

# References

1. Eberhard, K. The effects of visualization on judgment and decision-making: A systematic literature review. *Management Review Quarterly* (2021) doi:[10.1007/s11301-021-00235-8](https://doi.org/10.1007/s11301-021-00235-8).
2. Anscombe, F. J. **Graphs in Statistical Analysis**. *The American Statistician* **27**, 17–21 (1973).
3. Matejka, J. & Fitzmaurice, G. Same Stats, Different Graphs: Generating Datasets with Varied Appearance and Identical Statistics through Simulated Annealing. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* 1290–1294 (Association for Computing Machinery, 2017). doi:[10.1145/3025453.3025912](https://doi.org/10.1145/3025453.3025912).
4. Cairo, A. Download the Datasaurus: Never trust summary statistics alone; always visualize your data. (2016).
5. Snow, J. *On the mode of communication of cholera*. (John Churchill, 1855).
6. Nightingale, F. *Notes on Matters Affecting the Health, Efficiency and Hospital Administration of the British Army*. (Harrison & Sons, 1858).
7. Hans Rosling. The best stats you've ever seen [Video]. *The best stats you've ever seen* (2006).
8. Hawkins, E. Spiralling global temperatures | Climate Lab Book. (2016).
9. Pat Schloss. Recreating animated climate temperature spirals in R with Ggplot2 and gganimate (CC219). (2022).
10. Kosara, R. & Skau, D. Judgment Error in Pie Chart Variations. *EuroVis 2016 - Short Papers* 5 pages (2016) doi:[10.2312/EUROVISSHORT.20161167](https://doi.org/10.2312/EUROVISSHORT.20161167).
11. Cleveland, W. S. & McGill, R. **Graphical Perception: Theory, Experimentation, and Application to the Development of Graphical Methods**. *Journal of the American Statistical Association* **79**, 531–554 (1984).
12. Heer, J. & Bostock, M. Crowdsourcing graphical perception: Using mechanical turk to assess visualization design. In *Proceedings of the 28th international conference on Human factors in computing systems - CHI '10 203* (ACM Press, 2010). doi:[10.1145/1753326.1753357](https://doi.org/10.1145/1753326.1753357).
13. Beecham, R., Dykes, J., Hama, L. & Lomax, N. **On the Use of “Glyphmaps” for Analysing the Scale and Temporal Spread of COVID-19 Reported Cases**. *ISPRS International Journal of Geo-Information* **10**, 213 (2021).
14. Kosara, R. **More Than Meets the Eye: A Closer Look at Encodings in Visualization**. *IEEE Computer Graphics and Applications* **42**, 110–114 (2022).
15. Iker Rivas-González [@irg\_bio]. I am also joining the hexbin fever!  For this week's #TidyTuesday, I plotted the number of bee colonies in the US by year and season. It seems like cold and warm states have different patterns of seasonal changes. Code: [#RStats #DataViz #ggplot2 <https://t.co/OYgyg2azzM>. Twitter \(2022\).](https://github.com/rivasiker/TidyTuesday/blob/main/2022/2022-01-11/analysis_2022-01-11.Rmd)