# Cleveland's hierarchy



Cleveland's Graphical Features Hierarchy

**Best information transfer**

Position along a common scale: *bar chart, time series, scatter plot*

Position along nonaligned axis: *?????*

Length: *stacked bar charts, waterfall charts*

Angle / Slope: *pie charts*

Area: *pie charts, matrix charts, radar charts*

Volume: *3d charts*

Colour

**Worst information transfer**

Source: Presentation Graphics, Leland Wilkinson, SPSS Inc & Northwestern University
Revised 16Feb2010 td/wd

http://sfew.websitetoolbox.com/post/
clevelands-graphical-features-hierarchy-4598555

# Outline

# Scales

▶ The top of the hierarchy involves putting things on scales

# Scales

- The top of the hierarchy involves putting things on scales

- But what scale do we use?

# Scales

- The top of the hierarchy involves putting things on scales

- But what scale do we use?
  - Are our data anchored to zero?

# Scales

- The top of the hierarchy involves putting things on scales

- But what scale do we use?
  - Are our data anchored to zero?
    - If so, are we interested in differences or ratios?

# Scales

- The top of the hierarchy involves putting things on scales

- But what scale do we use?
    - Are our data anchored to zero?
        - If so, are we interested in differences or ratios?

    - Are they anchored somewhere else?

# Scales

- ▶ The top of the hierarchy involves putting things on scales

- ▶ But what scale do we use?
  - ▶ Are our data anchored to zero?
    - ▶ If so, are we interested in differences or ratios?
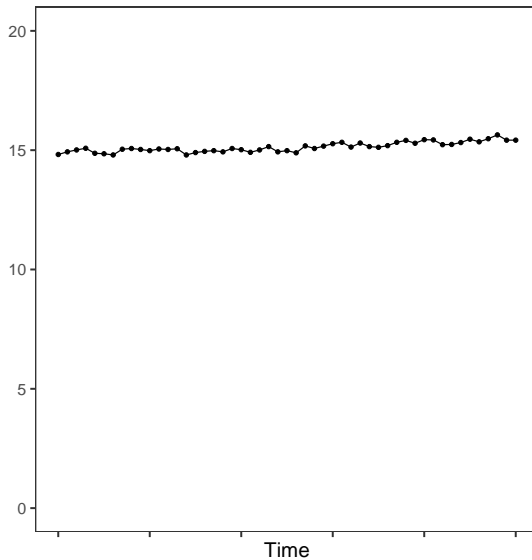
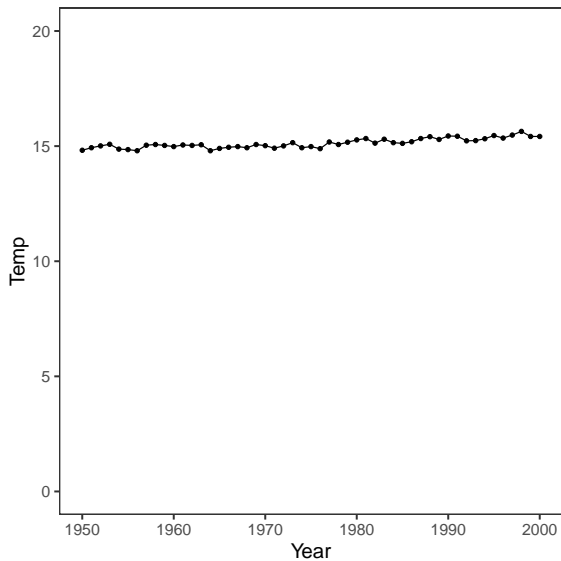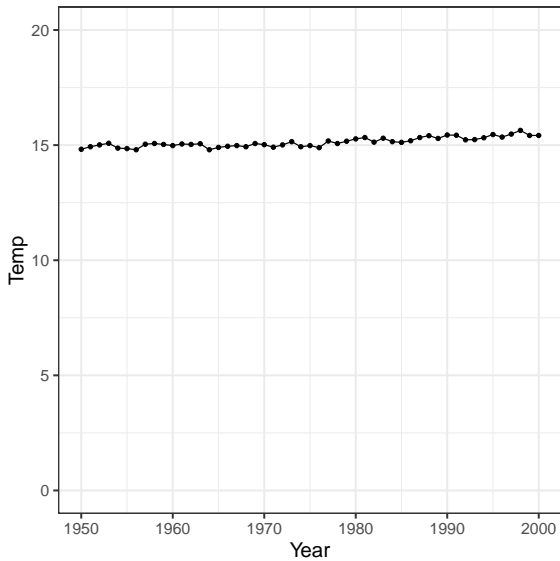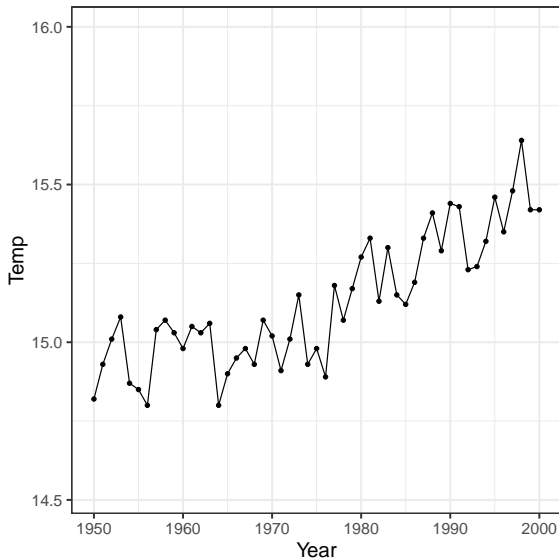  - ▶ Are they anchored somewhere else?
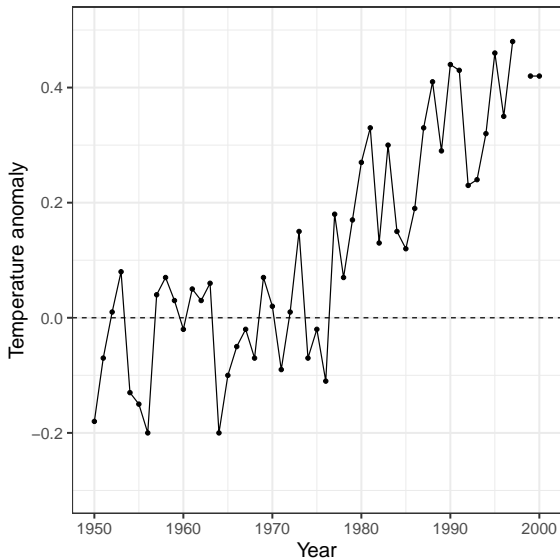
# Outline

# Golem bait call

# Global climate

# Global climate

# Global climate

# Global climate

# Climate lessons

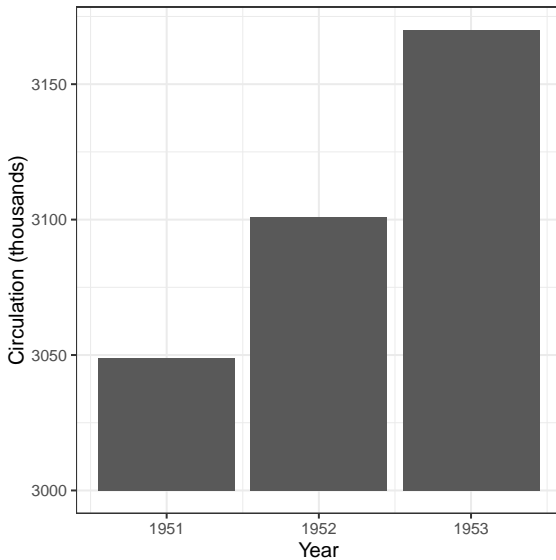- Choosing an anchor is a scientific decision

# Climate lessons

- Choosing an anchor is a scientific decision

- Remember: graphic design is communication

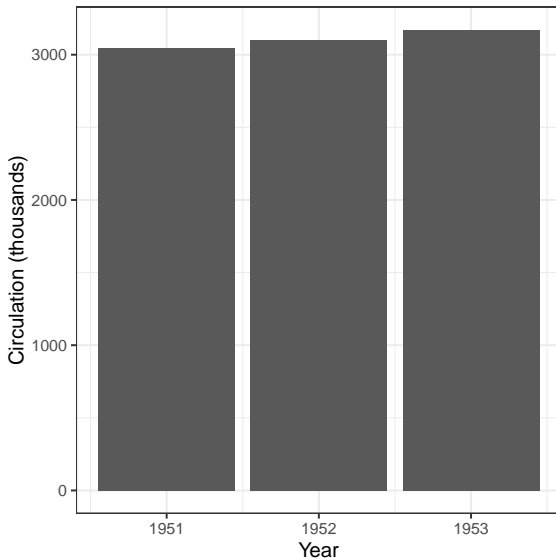# Climate lessons

► Choosing an anchor is a scientific decision

► Remember: graphic design is communication
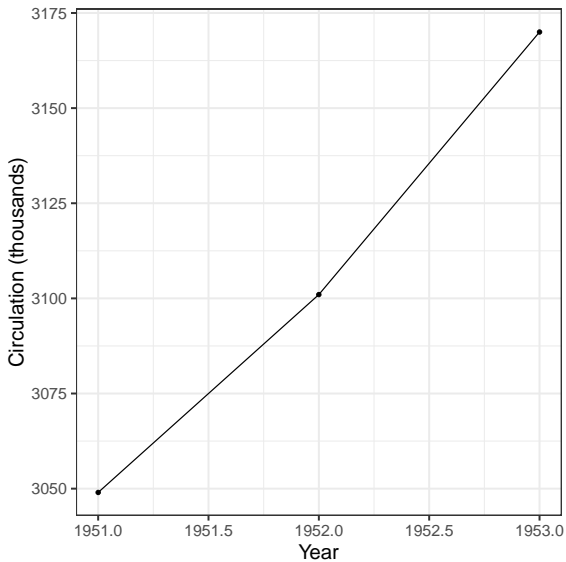
# Magazine circulation (advertisement)

# Magazine circulation (absolute amount)

# Magazine circulation (trend)

# Area and volume



STEEL CAPACITY ADDED

1930'S
10 MILLION TONS

1940'S
14¼ MILLION TONS

Adapted by courtesy of STEELWAYS.

*How to Lie with Statistics*

# Advertisement lessons

▶ Use area to indicate fair comparisons

# Advertisement lessons

- ▶ Use area to indicate fair comparisons
  - ▶ On a physical scale

# Advertisement lessons

- Use area to indicate fair comparisons
  - On a physical scale

- Areas that can be compared linearly should be preferred

# Advertisement lessons

- Use area to indicate fair comparisons
  - On a physical scale

- Areas that can be compared linearly should be preferred
  - Depends on importance of feature

# Advertisement lessons

- ▶ Use area to indicate fair comparisons
  - ▶ On a physical scale

- ▶ Areas that can be compared linearly should be preferred
  - ▶ Depends on importance of feature

- ▶ Avoid using (or hinting at) volume

# Advertisement lessons

- ▶ Use area to indicate fair comparisons
  - ▶ On a physical scale

- ▶ Areas that can be compared linearly should be preferred
  - ▶ Depends on importance of feature

- ▶ Avoid using (or hinting at) volume

# Outline

# Physical quantities

- ▶ 1 is to 10 as 10 is to what?

# Physical quantities

- 1 is to 10 as 10 is to what?
  - *

# Physical quantities

- 1 is to 10 as 10 is to what?
  - * If you said 19, you are thinking on a linear scale

# Physical quantities

- 1 is to 10 as 10 is to what?
  - \* If you said 19, you are thinking on a linear scale
  - \*

# Physical quantities

- 1 is to 10 as 10 is to what?
  - \* If you said 19, you are thinking on a linear scale
  - \* If you said 100, you are thinking on a log scale

# Physical quantities

- 1 is to 10 as 10 is to what?
  - \* If you said 19, you are thinking on a linear scale
  - \* If you said 100, you are thinking on a log scale

- The log scale is often good for physical quantities:

# Physical quantities

- 1 is to 10 as 10 is to what?
  - \* If you said 19, you are thinking on a linear scale
  - \* If you said 100, you are thinking on a log scale
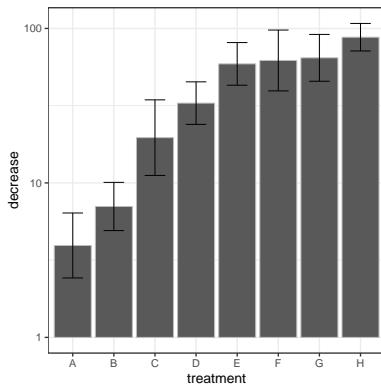
- The log scale is often good for physical quantities:
  - When zero means zero

# Physical quantities

- 1 is to 10 as 10 is to what?
  - * If you said 19, you are thinking on a linear scale
  - * If you said 100, you are thinking on a log scale
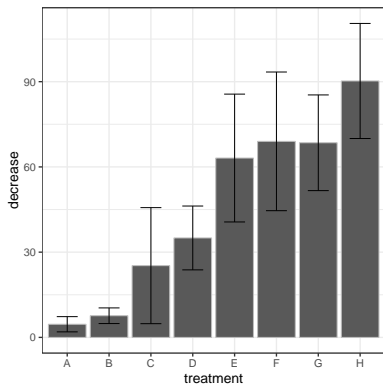
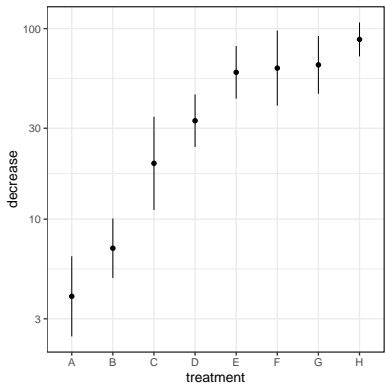- The log scale is often good for physical quantities:
  - When zero means zero

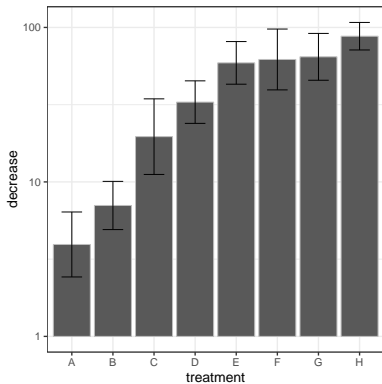# Log vs. linear

# Making room

# Data shape

▶ There are a lot of different ways to show data shape

# Data shape

- There are a lot of different ways to show data shape

- Choices will depend on your data set:

# Data shape

- There are a lot of different ways to show data shape

- Choices will depend on your data set:
  - Overall size

# Data shape

- There are a lot of different ways to show data shape

- Choices will depend on your data set:
  - Overall size
  - Number of replicates

# Data shape

- There are a lot of different ways to show data shape

- Choices will depend on your data set:
    - Overall size

    - Number of replicates
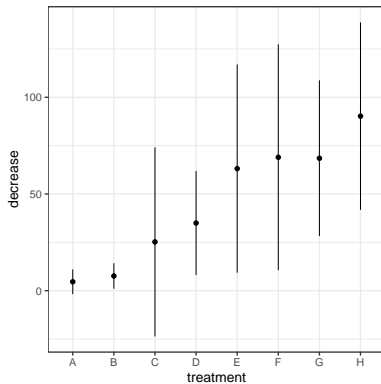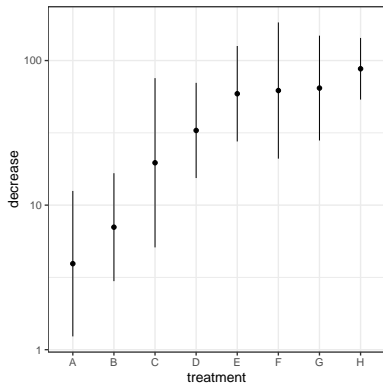
    - Number of levels, predictor variables, etc.

# Data shape

- There are a lot of different ways to show data shape

- Choices will depend on your data set:
  - Overall size
  - Number of replicates
  - Number of levels, predictor variables, etc.
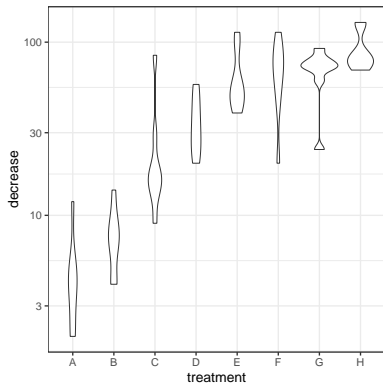
# Different scales

# More detail

# Orchard lessons

▶ Choices about log vs. linear scale are scientific choices

# Orchard lessons

- Choices about log vs. linear scale are scientific choices
  - Neither is more valid, or closer to the data

# Orchard lessons

- Choices about log vs. linear scale are scientific choices
  - Neither is more valid, or closer to the data

- You can also make choices about

# Orchard lessons

- Choices about log vs. linear scale are scientific choices
  - Neither is more valid, or closer to the data

- You can also make choices about
  - sending a simple message

# Orchard lessons

- Choices about log vs. linear scale are scientific choices
  - Neither is more valid, or closer to the data

- You can also make choices about
  - sending a simple message
  - providing more information about shape

# Orchard lessons

- Choices about log vs. linear scale are scientific choices
    - Neither is more valid, or closer to the data

- You can also make choices about
    - sending a simple message
    - providing more information about shape

- Log scales are almost never physical

# Orchard lessons

- Choices about log vs. linear scale are scientific choices
  - Neither is more valid, or closer to the data

- You can also make choices about
  - sending a simple message
  - providing more information about shape

- Log scales are almost never physical
  - Don't mislead with area information on a log scale

# Orchard lessons

- Choices about log vs. linear scale are scientific choices
  - Neither is more valid, or closer to the data

- You can also make choices about
  - sending a simple message
  - providing more information about shape

- Log scales are almost never physical
  - Don't mislead with area information on a log scale

# Probabilities

▶ 1% is to 2% as 50% is to what?

# Probabilities

- 1% is to 2% as 50% is to what?
  - *

# Probabilities

- 1% is to 2% as 50% is to what?
  - \* 51% is way too small

# Probabilities

- 1% is to 2% as 50% is to what?
  - * 51% is way too small
  - *

# Probabilities

- 1% is to 2% as 50% is to what?
  - \* 51% is way too small
  - \* 100% is way too large

# Probabilities

- 1% is to 2% as 50% is to what?
  - * 51% is way too small
  - * 100% is way too large

- The natural distance to use on a probability scale is log odds

# Probabilities

- 1% is to 2% as 50% is to what?
  - * 51% is way too small
  - * 100% is way too large

- The natural distance to use on a probability scale is log odds
  - *

# Probabilities

- 1% is to 2% as 50% is to what?
  - * 51% is way too small
  - * 100% is way too large

- The natural distance to use on a probability scale is log odds
  - * 1% is to 2% as 50% is to 67%

# Probabilities

- 1% is to 2% as 50% is to what?
  - * 51% is way too small
  - * 100% is way too large

- The natural distance to use on a probability scale is log odds
  - * 1% is to 2% as 50% is to 67%
  - *

# Probabilities

- 1% is to 2% as 50% is to what?
  - \* 51% is way too small
  - \* 100% is way too large

- The natural distance to use on a probability scale is log odds
  - \* 1% is to 2% as 50% is to 67%
  - \* . . . as 2% is to 4%

# Probabilities

- 1% is to 2% as 50% is to what?
  - * 51% is way too small
  - * 100% is way too large

- The natural distance to use on a probability scale is log odds
  - * 1% is to 2% as 50% is to 67%
  - * . . . as 2% is to 4%
  - *

# Probabilities

- 1% is to 2% as 50% is to what?
  - * 51% is way too small
  - * 100% is way too large

- The natural distance to use on a probability scale is log odds
  - * 1% is to 2% as 50% is to 67%
  - * ... as 2% is to 4%
  - * ... as 98% is to 99%

# Probabilities

- 1% is to 2% as 50% is to what?
  - \* 51% is way too small
  - \* 100% is way too large

- The natural distance to use on a probability scale is log odds
  - \* 1% is to 2% as 50% is to 67%
  - \* . . . as 2% is to 4%
  - \* . . . as 98% is to 99%

# Odds

▶ Odds are a ratio between the probability of something and the probability of its opposite:

# Odds

- Odds are a ratio between the probability of something and the probability of its opposite:
  - $o = p/(1 - p)$

# Odds

- Odds are a ratio between the probability of something and the probability of its opposite:
  - $o = p/(1 - p)$

- Log odds give a natural distance on probability space

# Odds

- Odds are a ratio between the probability of something and the probability of its opposite:
  - $o = p/(1 - p)$

- Log odds give a natural distance on probability space

# Extreme values

▶ Our transformations take extreme values to infinity.

# Extreme values

- ▶ Our transformations take extreme values to infinity.

- ▶ Use link functions

# Extreme values

- Our transformations take extreme values to infinity.

- Use link functions
  - this is like using estimated values instead of observed

# Extreme values

- ▶ Our transformations take extreme values to infinity.

- ▶ Use link functions
  - ▶ this is like using estimated values instead of observed
    - ▶ rarely infinite

# Extreme values

- Our transformations take extreme values to infinity.

- Use link functions
  - this is like using estimated values instead of observed
    - rarely infinite
    - matches analysis

# Extreme values

- Our transformations take extreme values to infinity.

- Use link functions
  - this is like using estimated values instead of observed
    - rarely infinite
    - matches analysis

- Extend the scale (e.g., use $\log(1 + x)$ instead of $\log(x)$)

# Extreme values

- Our transformations take extreme values to infinity.

- Use link functions
  - this is like using estimated values instead of observed
    - rarely infinite
    - matches analysis

- Extend the scale (e.g., use $\log(1 + x)$ instead of $\log(x)$)
  - This usually involves arbitrary choices

# Extreme values

- Our transformations take extreme values to infinity.

- Use link functions
  - this is like using estimated values instead of observed
    - rarely infinite
    - matches analysis

- Extend the scale (e.g., use $\log(1 + x)$ instead of $\log(x)$)
  - This usually involves arbitrary choices
  - Should often be *avoided* for analysis

# Extreme values

- Our transformations take extreme values to infinity.

- Use link functions
  - this is like using estimated values instead of observed
    - rarely infinite
    - matches analysis

- Extend the scale (e.g., use $\log(1 + x)$ instead of $\log(x)$)
  - This usually involves arbitrary choices
  - Should often be *avoided* for analysis
  - But can be good for visualization

# Extreme values

- ▶ Our transformations take extreme values to infinity.

- ▶ Use link functions
  - ▶ this is like using estimated values instead of observed
    - ▶ rarely infinite
    - ▶ matches analysis

- ▶ Extend the scale (e.g., use $\log(1 + x)$ instead of $\log(x)$)
  - ▶ This usually involves arbitrary choices
  - ▶ Should often be *avoided* for analysis
  - ▶ But can be good for visualization

# Outline

# Rote analysis vs. snooping

# Spurious correlations

There's a whole website about this

# What can you do?

The best you can

▶ Identify scientific questions

# What can you do?

The best you can

- ▶ Identify scientific questions

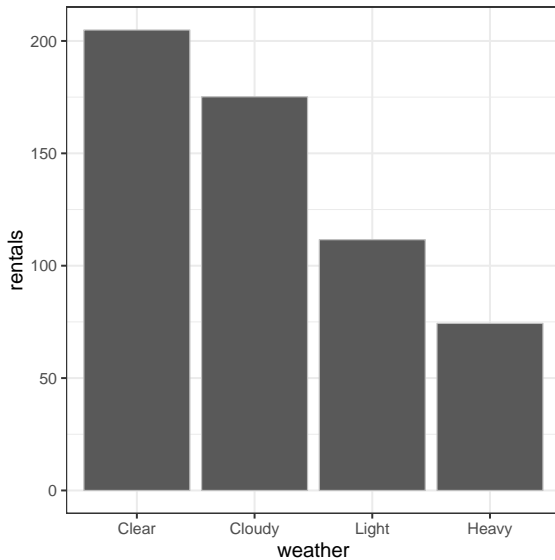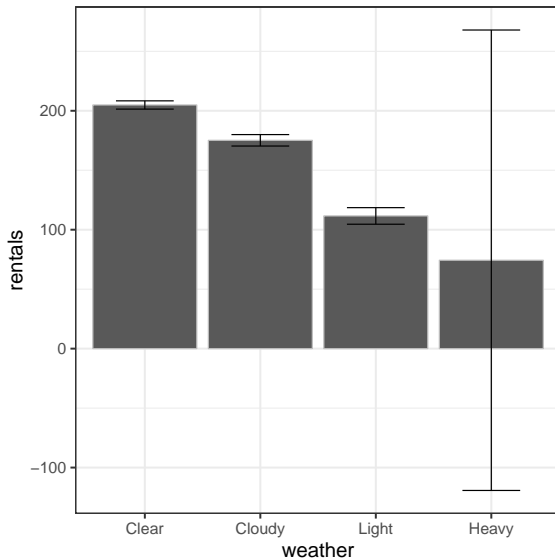- ▶ Distinguish between exploratory and confirmatory analysis

# What can you do?
The best you can

- Identify scientific questions

- Distinguish between exploratory and confirmatory analysis

- Pre-register studies when possible

# What can you do?
## The best you can

- ▶ Identify scientific questions

- ▶ Distinguish between exploratory and confirmatory analysis

- ▶ Pre-register studies when possible

- ▶ Keep an exploration and analysis journal

# What can you do?
## The best you can

- ▶ Identify scientific questions

- ▶ Distinguish between exploratory and confirmatory analysis

- ▶ Pre-register studies when possible

- ▶ Keep an exploration and analysis journal

- ▶ Explore predictors and responses separately at first

# What can you do?
## The best you can
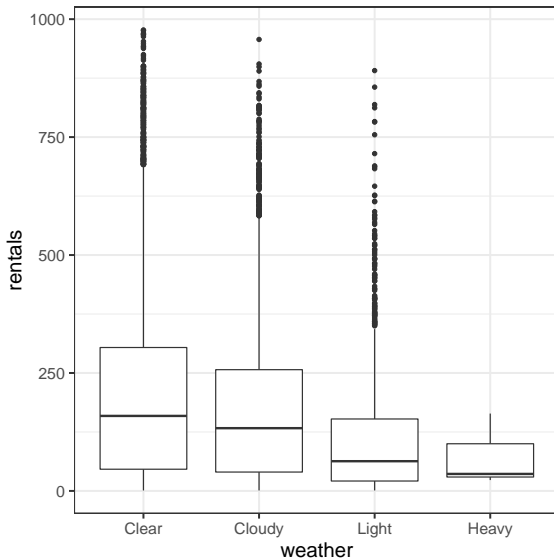
- ▶ Identify scientific questions

- ▶ Distinguish between exploratory and confirmatory analysis

- ▶ Pre-register studies when possible

- ▶ Keep an exploration and analysis journal

- ▶ Explore predictors and responses separately at first
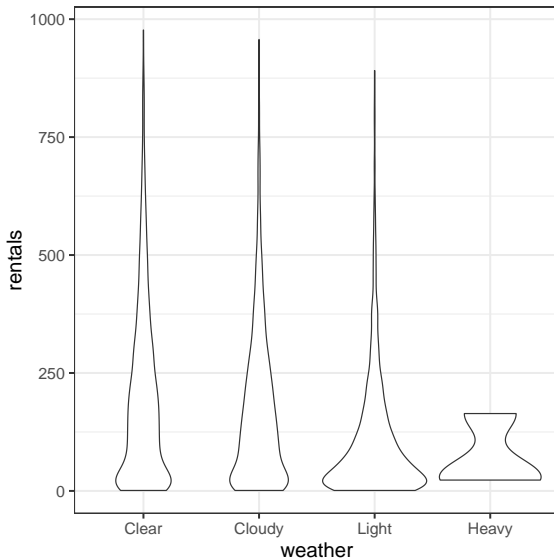
# Bike example

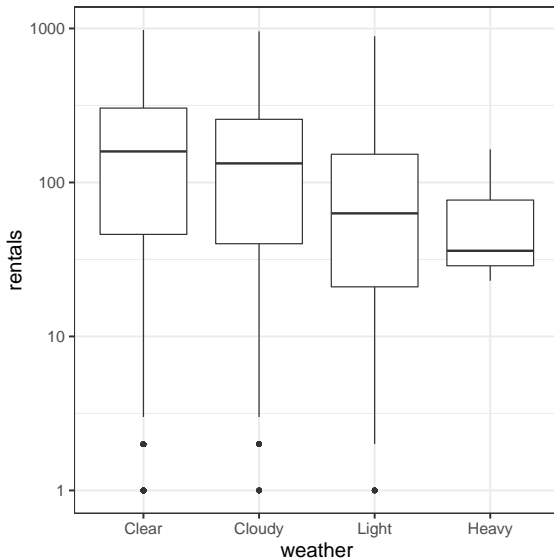# Standard errors

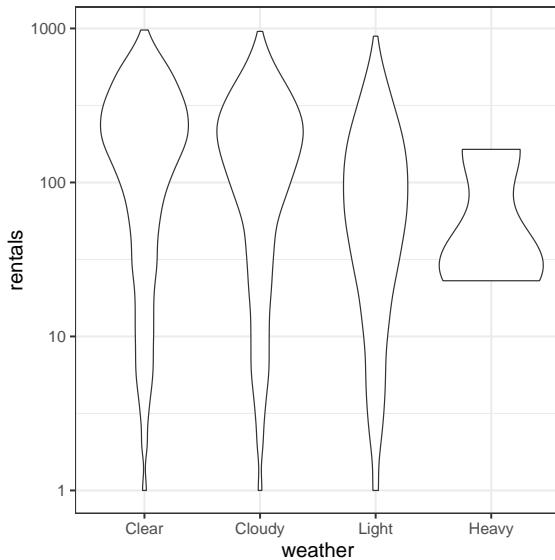# Standard errors
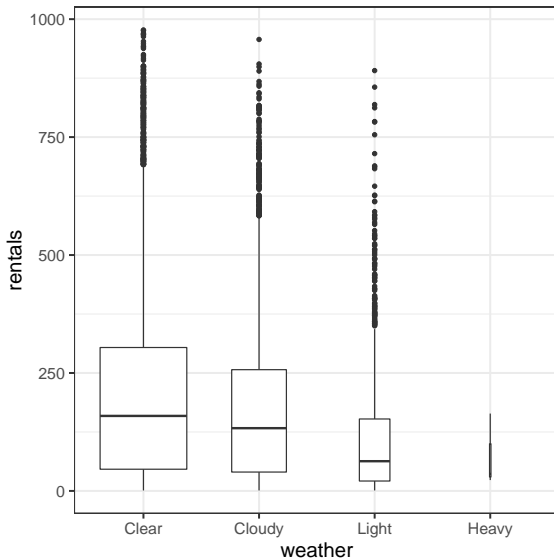
# Standard deviations

# Data shape

# Data shape

# Data shape

# Data shape

# Data shape and weight

# Log scales

▶ In general:

# Log scales

- In general:
  - If your logged data span $< 3$ decades, use human-readable numbers (e.g., 10-5000 kilotons per hectare)

# Log scales

- In general:
  - If your logged data span $< 3$ decades, use human-readable numbers (e.g., 10-5000 kilotons per hectare)
  - If not, just embrace "logs" (log10 particles per ul is from 3–8)

# Log scales

- In general:
  - If your logged data span $< 3$ decades, use human-readable numbers (e.g., 10-5000 kilotons per hectare)
  - If not, just embrace "logs" (log10 particles per ul is from 3–8)
    - But remember these are not physical values

# Log scales

- In general:
  - If your logged data span $< 3$ decades, use human-readable numbers (e.g., 10-5000 kilotons per hectare)
  - If not, just embrace "logs" (log10 particles per ul is from 3–8)
    - But remember these are not physical values

- I love natural logs, but not as axis values

# Log scales

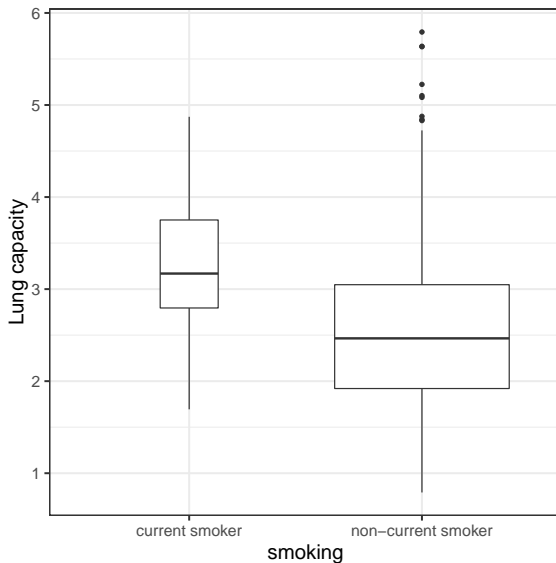- In general:
    - If your logged data span $< 3$ decades, use human-readable numbers (e.g., 10-5000 kilotons per hectare)
    - If not, just embrace "logs" (log10 particles per ul is from 3–8)
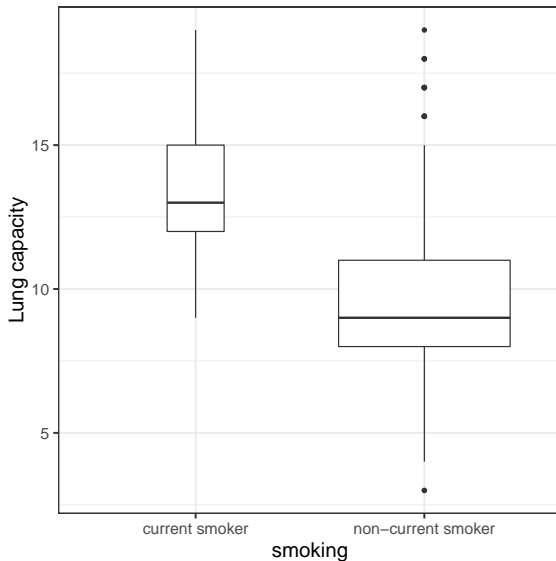        - But remember these are not physical values

- I love natural logs, but not as axis values

# Outline

# Smoking data

# Smoking data

# Scatter plots

▶ Depending on how many data points you have, scatter plots may indicate relationships clearly

# Scatter plots

- Depending on how many data points you have, scatter plots may indicate relationships clearly

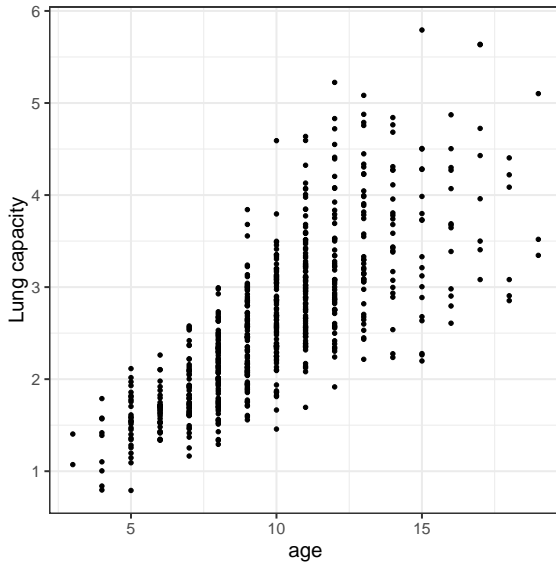- They can often be improved with trend interpolations

# Scatter plots

▶ Depending on how many data points you have, scatter plots may indicate relationships clearly

▶ They can often be improved with trend interpolations
  ▶ Interpolations may be particularly good for discrete responses (count or true-false)
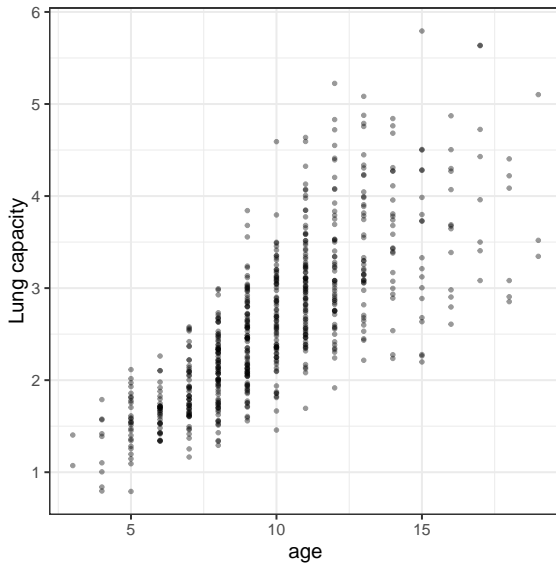
# Scatter plots

- Depending on how many data points you have, scatter plots may indicate relationships clearly

- They can often be improved with trend interpolations
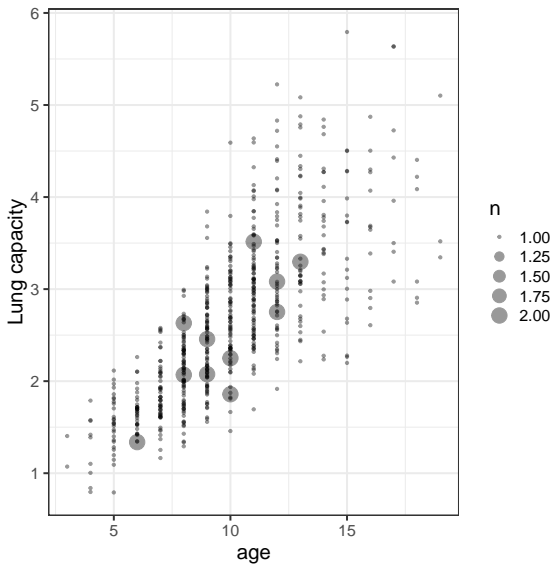  - Interpolations may be particularly good for discrete responses (count or true-false)
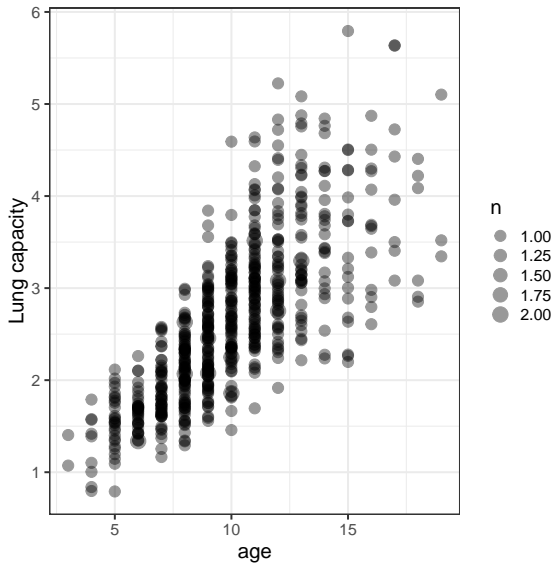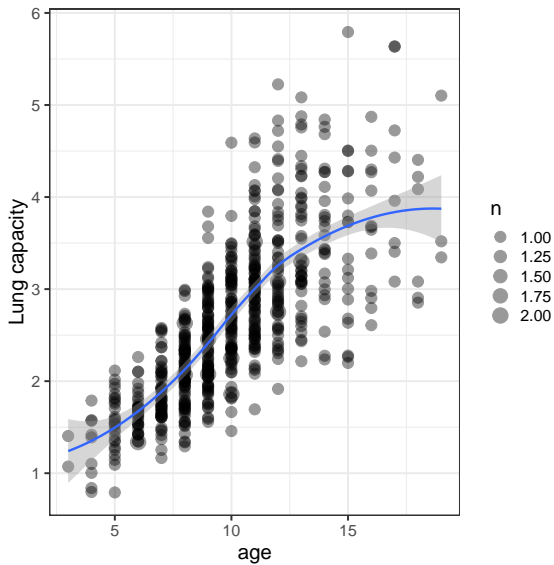
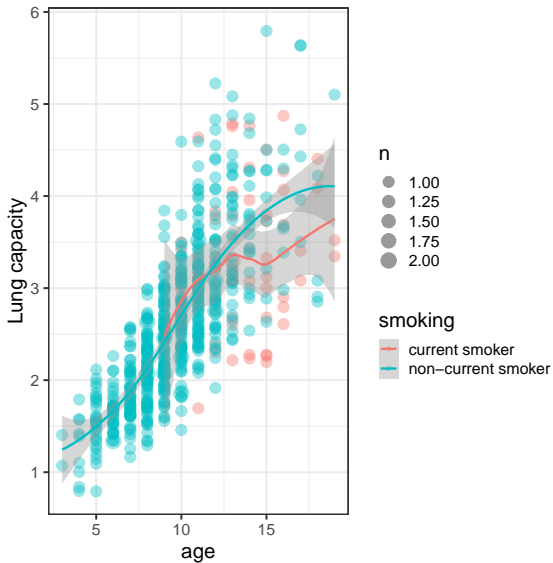# Scatter plot

# Seeing the density better

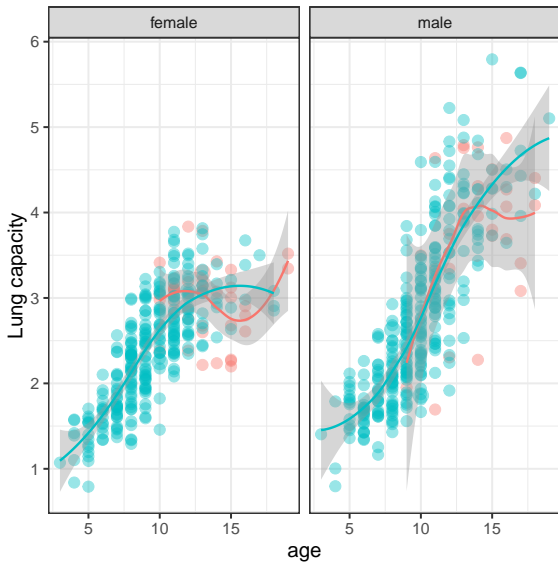# Seeing the density worse

# Use area in a principled way!

# A loess trend line

# Two loess trend lines

# Many loess trend lines

# Density plots

▶ Contours

# Density plots

- Contours
  - use `_density_2d()` to fit a two-dimensional kernel to the density

# Density plots

- Contours
  - use `_density_2d()` to fit a two-dimensional kernel to the density

- hexes

# Density plots

- Contours
  - use `_density_2d()` to fit a two-dimensional kernel to the density

- hexes
  - use geom_hex to plot densities using hexes

# Density plots

- Contours
  - use `_density_2d()` to fit a two-dimensional kernel to the density

- hexes
  - use `geom_hex` to plot densities using hexes
  - this can also be done using rectangles for data with more discrete values

# Density plots
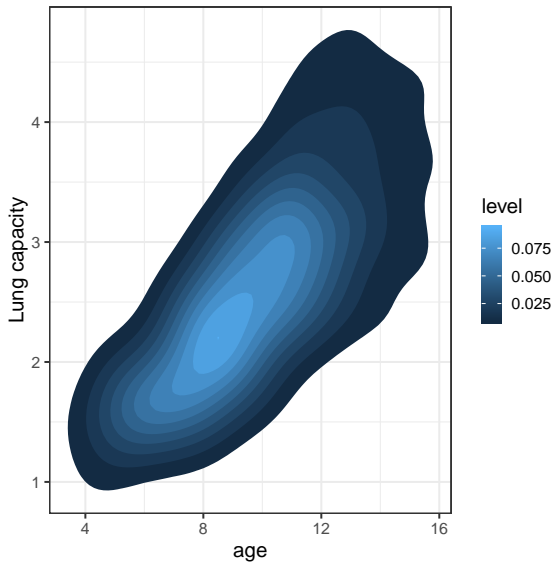
- Contours
  - use `_density_2d()` to fit a two-dimensional kernel to the density

- hexes
  - use `geom_hex` to plot densities using hexes

  - this can also be done using rectangles for data with more discrete values

# Contours

# Contours

# Hexes

# Hexes

# Color principles

- Use clear gradients

# Color principles

- ► Use clear gradients

- ► If zero has a physical meaning (like density), go in just one direction

# Color principles

- Use clear gradients

- If zero has a physical meaning (like density), go in just one direction
  - e.g., white to blue, white to red

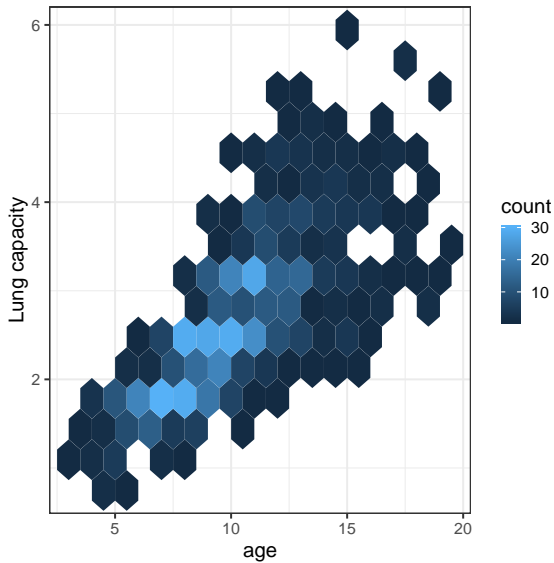# Color principles

- ▶ Use clear gradients

- ▶ If zero has a physical meaning (like density), go in just one direction
    - ▶ e.g., white to blue, white to red
    - ▶ If the map contrasts with a background, zero should match the background

# Color principles

- ▶ Use clear gradients

- ▶ If zero has a physical meaning (like density), go in just one direction
  - ▶ e.g., white to blue, white to red
  - ▶ If the map contrasts with a background, zero should match the background

- ▶ If there's a natural *middle*, you can use blue to white to red, or something similar

# Color principles

- ▶ Use clear gradients

- ▶ If zero has a physical meaning (like density), go in just one direction
  - ▶ e.g., white to blue, white to red
  - ▶ If the map contrasts with a background, zero should match the background

- ▶ If there's a natural *middle*, you can use blue to white to red, or something similar

# Principles

▶ Graphs tell stories better than tables do

# Principles

▶ Graphs tell stories better than tables do
  ▶ Use graphs to illustrate comparisons

# Principles

- Graphs tell stories better than tables do
  - Use graphs to illustrate comparisons
  - Be careful about *units*

# Principles

- Graphs tell stories better than tables do
  - Use graphs to illustrate comparisons
  - Be careful about *units*

- Distinguish between (scientific) variables and (statistical) parameters

# Principles

- Graphs tell stories better than tables do
  - Use graphs to illustrate comparisons
  - Be careful about *units*

- Distinguish between (scientific) variables and (statistical) parameters

- Show data when you can do it without obscuring the key patterns

# Principles

- ► Graphs tell stories better than tables do
  - ► Use graphs to illustrate comparisons
  - ► Be careful about *units*

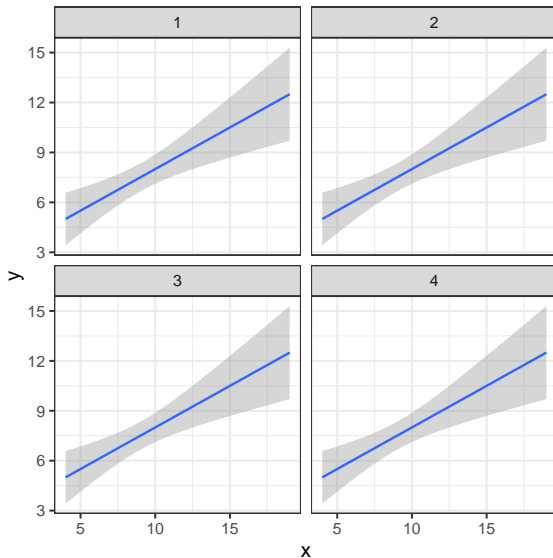- ► Distinguish between (scientific) variables and (statistical) parameters

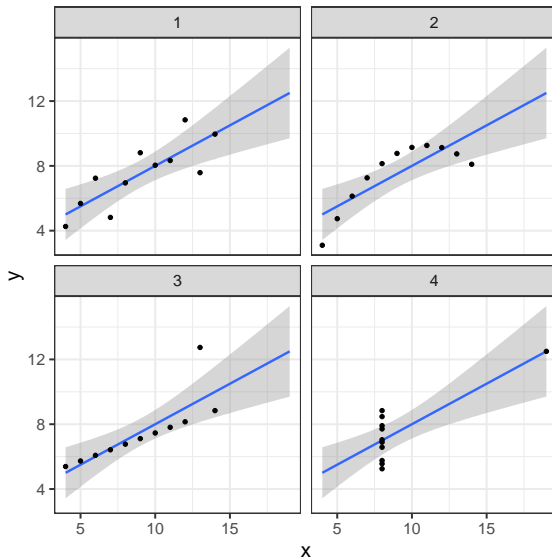- ► Show data when you can do it without obscuring the key patterns

# Choosing what to show

# Choosing what to show

# Avoiding choices

- Provide users with alternatives

# Avoiding choices

- ▶ Provide users with alternatives
  - ▶ Supplementary material for the curious

# Avoiding choices

- Provide users with alternatives
  - Supplementary material for the curious

- Avoid choices by providing more information

# Avoiding choices

- ▶ Provide users with alternatives
  - ▶ Supplementary material for the curious

- ▶ Avoid choices by providing more information
  - ▶ Use more than one figure

# Avoiding choices

- ▶ Provide users with alternatives
    - ▶ Supplementary material for the curious

- ▶ Avoid choices by providing more information
    - ▶ Use more than one figure
    - ▶ Or dynamic features in figures

# Avoiding choices

- ▶ Provide users with alternatives
  - ▶ Supplementary material for the curious

- ▶ Avoid choices by providing more information
  - ▶ Use more than one figure
  - ▶ Or dynamic features in figures

# Conclusions

- Give thought to your goals

# Conclusions

▶ Give thought to your goals

▶ Give thought to your decisions

# Conclusions

▶ Give thought to your goals

▶ Give thought to your decisions

▶ Be conscious when you are withholding information

# Conclusions

▶ Give thought to your goals

▶ Give thought to your decisions

▶ Be conscious when you are withholding information
  ▶ Be willing to use more than one picture

# Conclusions

- ▶ Give thought to your goals

- ▶ Give thought to your decisions

- ▶ Be conscious when you are withholding information
  - ▶ Be willing to use more than one picture
  - ▶ Use dynamic features to give access to detail

# Conclusions

▶ Give thought to your goals

▶ Give thought to your decisions

▶ Be conscious when you are withholding information
  ▶ Be willing to use more than one picture
  ▶ Use dynamic features to give access to detail