

# [IGP 484] Quantitative Biology: Data Analysis for Life Scientists

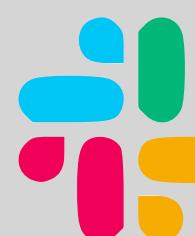
Dr. Katie Evans



Tu/Th 1:30-3:00, Hughes Auditorium



[github.com/katieseavans/IGP biostatistics](https://github.com/katieseavans/IGP_biotatistics)



biostats-484



Tu/Th 1:30-3:00, Hughes Auditorium\*



[github.com/katieselevans/IGP biostatistics](https://github.com/katieselevans/IGP_biolab)

- Weekly homework assignments (40%)
- Class participation (5%)
- Midterm exam (15%)
- Final exam (15%)
- Final project (25%)

\*9/23 and 11/11 – Class will meet in Daniel Hale Williams, McGaw 2-320



Tu/Th 1:30-3:00, Hughes Auditorium



[github.com/katieselevans/IGP biostatistics](https://github.com/katieselevans/IGP_biolab)

- **Homework assignments**

- Must be completed using R
- Make sure to show ALL your work for credit
- Submitted electronically as PDF or markdown

- **Final project**

- Choose a dataset (rotation, internet, old publication...)
- Perform a statistical analysis
- Write up a report with conclusions
- *More to come in second half of class, but start thinking about ideas!*

# Meet the class TA's!



**Saya**



**Jiexi**



**Sam**



**Yidan**



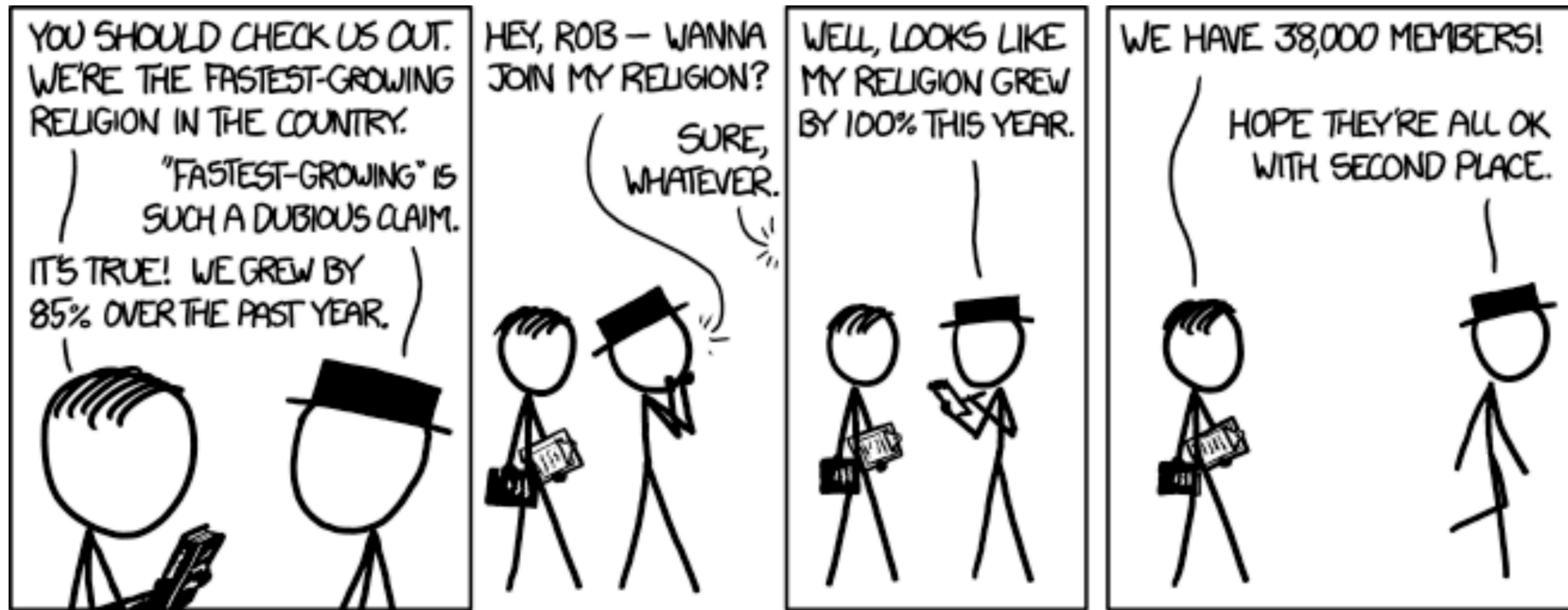
**TA office hours:** Tuesday 3-4 (TBD); Friday 1:30-2:30 (Zoom)

*(Or by appointment)*

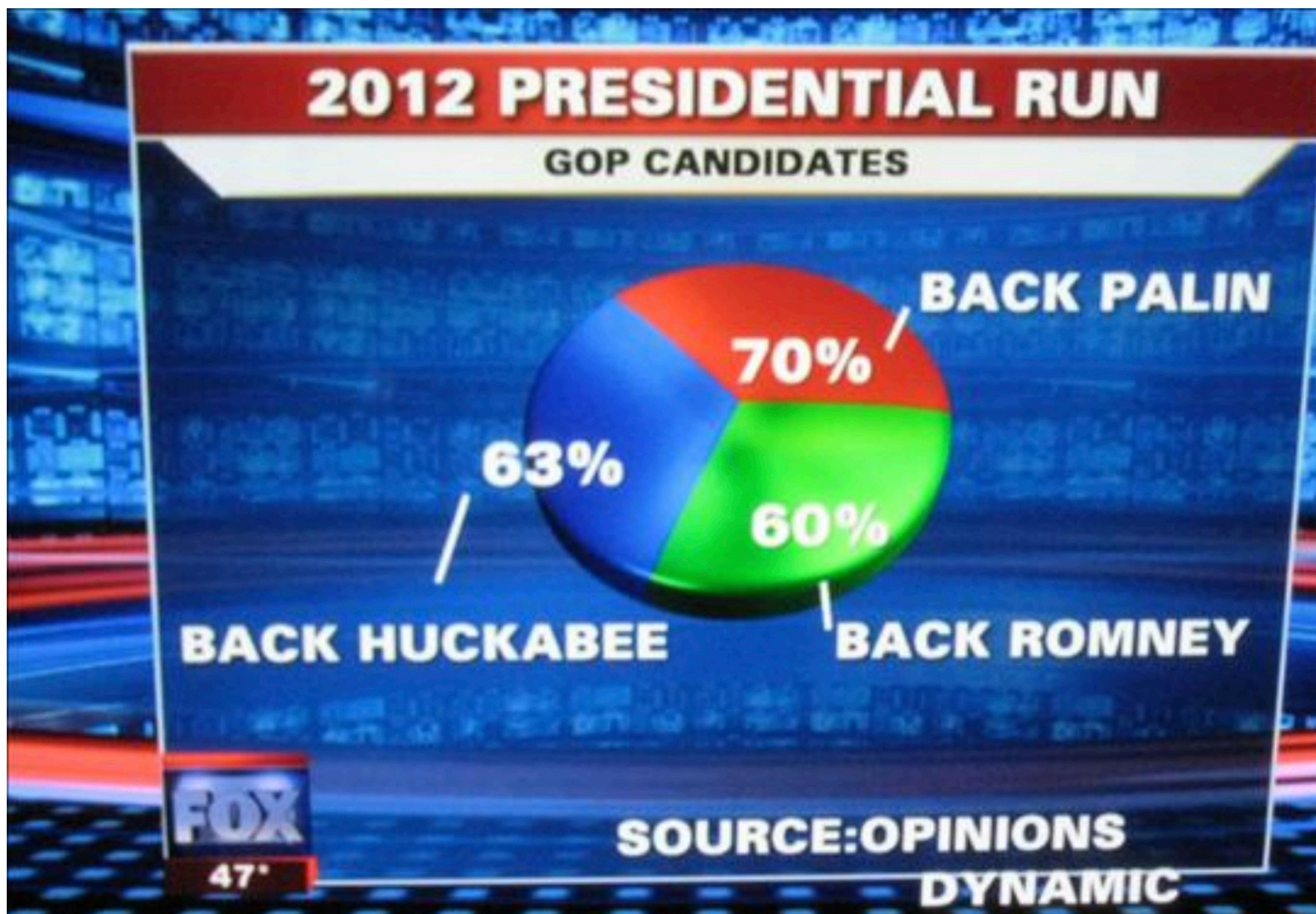
Date	Class	Lecture	Topics	Textbook Sections
September 21, 2021	Lecture	<a href="#">Lecture 1: Introduction</a>	Class intro, what is "statistics"?, displaying and summarizing data, getting started with R	1.1-1.3, 2.1-2.7
September 23, 2021	Lecture	<a href="#">Lecture 2: Probability</a>	notation, independence, conditional probability	3.1-3.3
September 28, 2021	Lecture	<a href="#">Lecture 3: Probability (part 2)</a>	conditional probability, Bayes theorem, practical implications	3.3-3.5
September 30, 2021	Lecture	<a href="#">Lecture 4: Distributions</a>	binomial and normal	3.6, 4.1-4.3
October 5, 2021	Practicum	<a href="#">Practicum 1: Learning R</a>	Data tidying with tidyverse, plotting with ggplot2, thinking like a data scientist, Rmarkdown	NA
October 7, 2021	Lecture	<a href="#">Lecture 5: Distributions (part 2)</a>	Central Limit Theorem, sampling	5.1-5.4
October 12, 2021	Lecture	<a href="#">Lecture 6: Estimation, testing, and p-values</a>	estimation, testing, p-values	4.4, 6.1-6.5
October 14, 2021	Lecture	<a href="#">Lecture 7: One-sample t-tests</a>	estimating and testing means	6.6-6.7, 7.2-7.6, 7.8
October 19, 2021	Lecture	<a href="#">Lecture 8: Two-sample comparisons</a>	unpaired t-tests, paired t-tests, power	7.2-7.6, 8.2-8.3, 7.7
October 21, 2021	Lecture	<a href="#">Lecture 9: Nonparametric alternatives</a>	permutation, Wilcoxon/Mann-Whitney, signed rank	7.1, 7.10, 8.4-8.5
October 26, 2021	Review	<a href="#">Q&amp;A: Project discussion and Midterm review</a>	Bring your questions!	Ch. 1-8
October 28, 2021	Exam	<a href="#">Midterm exam</a>	NA	NA

November 2, 2021	Lecture	<a href="#">Lecture 10: Multiple hypotheses</a>	Urn models, "enrichment", Fischer's exact test	10.4
November 4, 2021		<a href="#">Practicum 2: Analyzing gene expression data with R</a>	NA	NA
November 9, 2021	Lecture	<a href="#">Lecture 11: Categorical data</a>	contingency tables, chi-squared tests, relative risks and odds ratio	9.1-9.2, 9.4, 10.1-10.3, 10.5-10.6, 10.9-10.10
November 11, 2021	Lecture	<a href="#">Lecture 12: Relationships in data</a>	independence, covariance, and correlation	12.1-12.2
November 16, 2021	Lecture	<a href="#">Lecture 13: Regression models</a>	linear assumptions, interpretation, limitations, transformations	12.4-12.6
November 18, 2021	Lecture	<a href="#">Lecture 14: Models with categorical covariates</a>	indicator/"dummy" variables	11.1-11.5
November 23, 2021	Lecture	<a href="#">Lecture 15: ANOVA</a>	one-way ANOVA, two-way ANOVA	11.7-11.8
November 25, 2021		No class: Thanksgiving!	NA	NA
November 30, 2021	Practicum	<a href="#">Practicum 3: Model fitting in R</a>	NA	NA
December 2, 2021	Review	<a href="#">Q&amp;A: "Overflow" topics and final review</a>	Cumulative!	NA
December 7, 2021	Exam	<a href="#">Final Exam</a>	Cumulative!	NA

# Misleading statistics

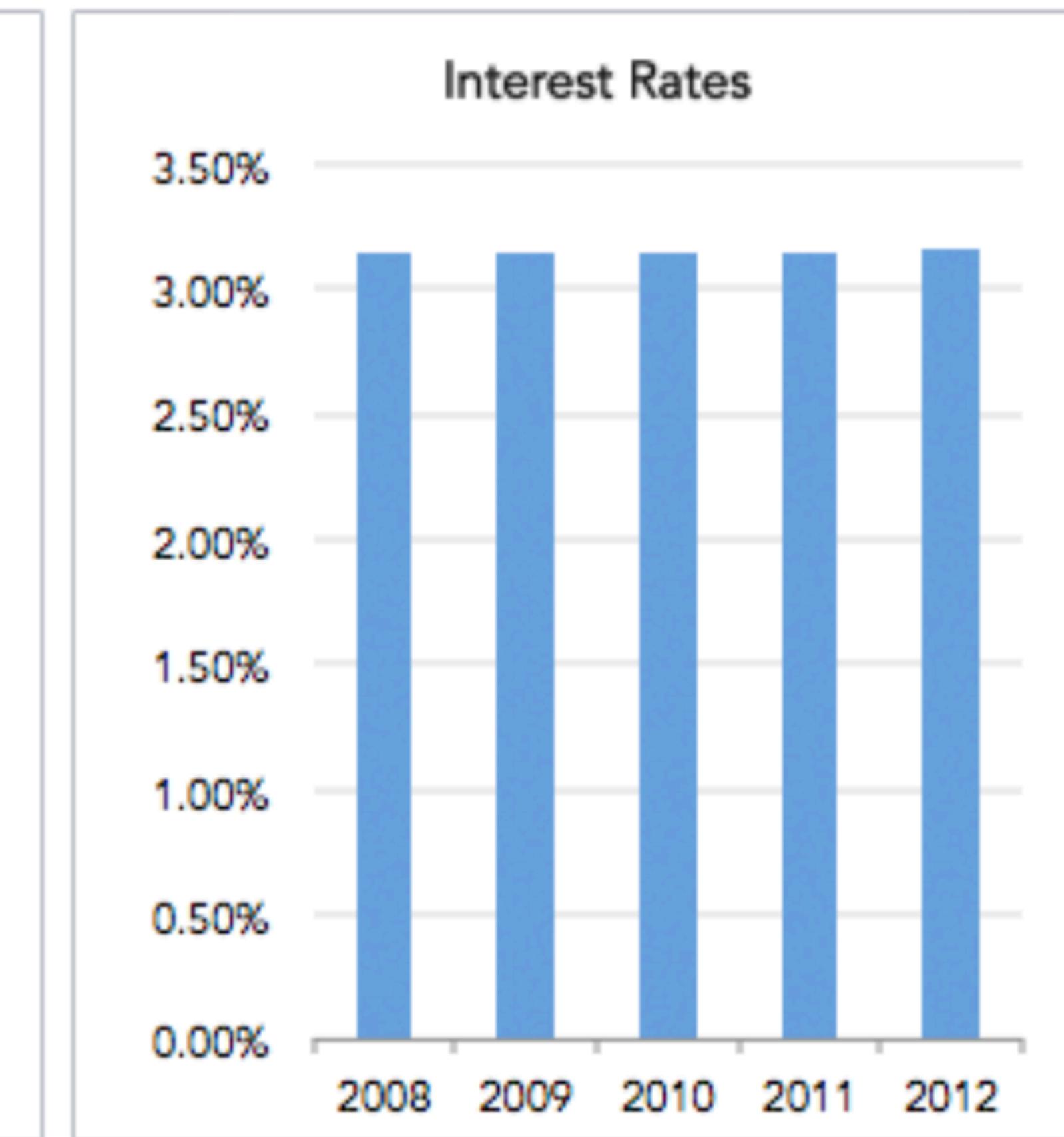
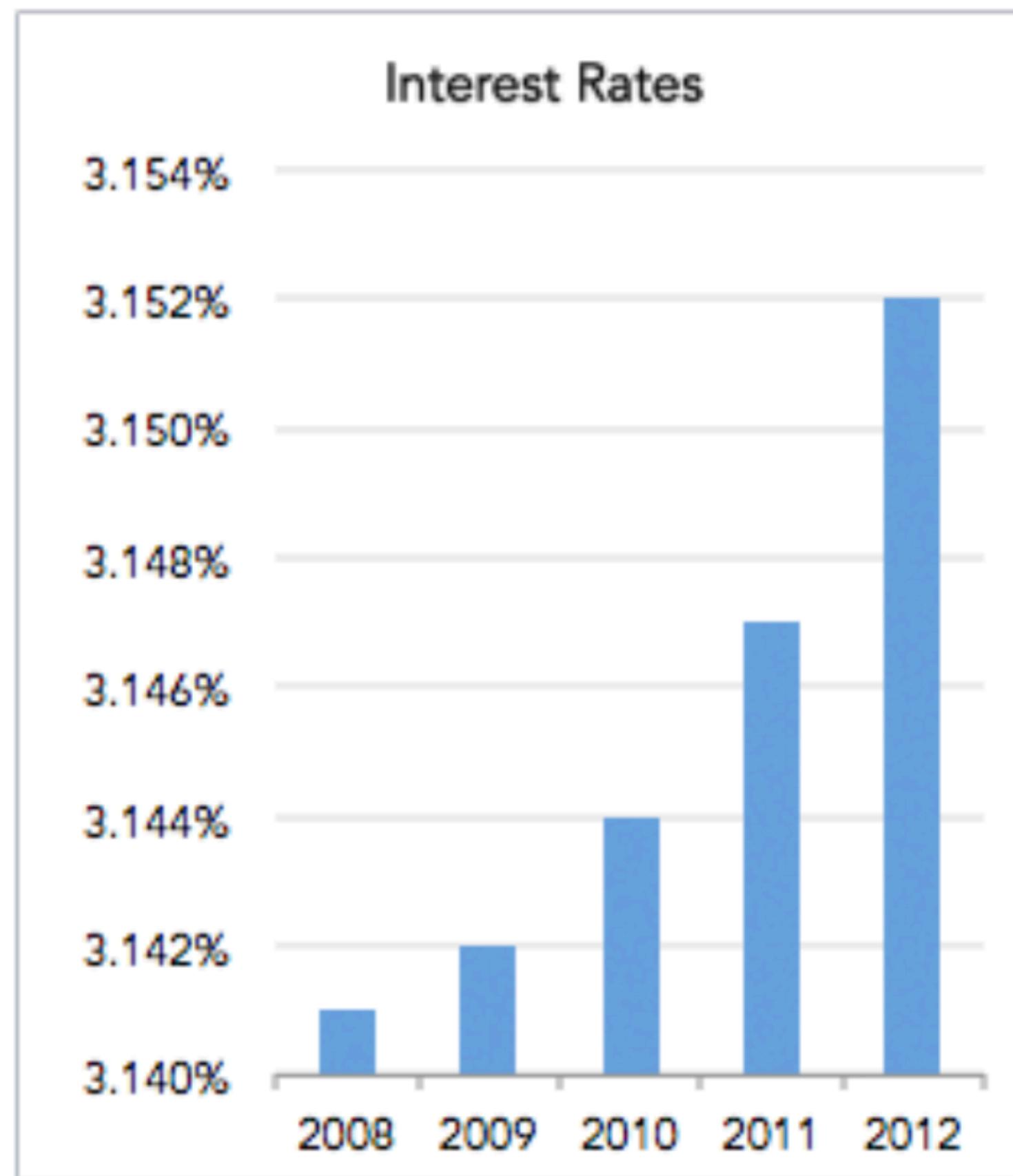


# Misleading statistics



# Misleading statistics

**Same Data, Different Y-Axis**



# Misleading statistics



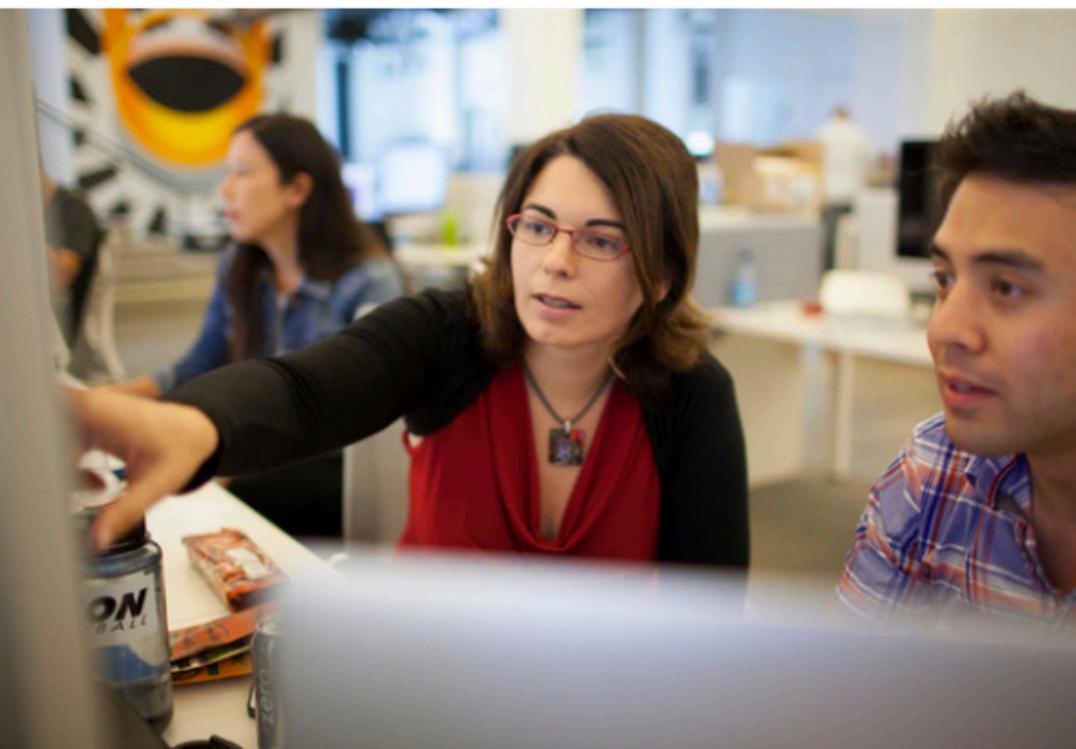
# Reproducible research



The New York Times

**For Big-Data Scientists, ‘Janitor Work’ Is Key Hurdle to Insights**

By STEVE LOHR AUG. 17, 2014

A photograph showing two data scientists, Monica Rogati and Brian Wilt, sitting at a desk in an office environment. Monica is wearing glasses and a red top, while Brian is wearing a plaid shirt. They appear to be engaged in a discussion or analysis.

Monica Rogati, Jawbone's vice president for data science, with Brian Wilt, a senior data scientist.  
Peter DaSilva for The New York Times



FORTUNE

## Big data's dirty problem

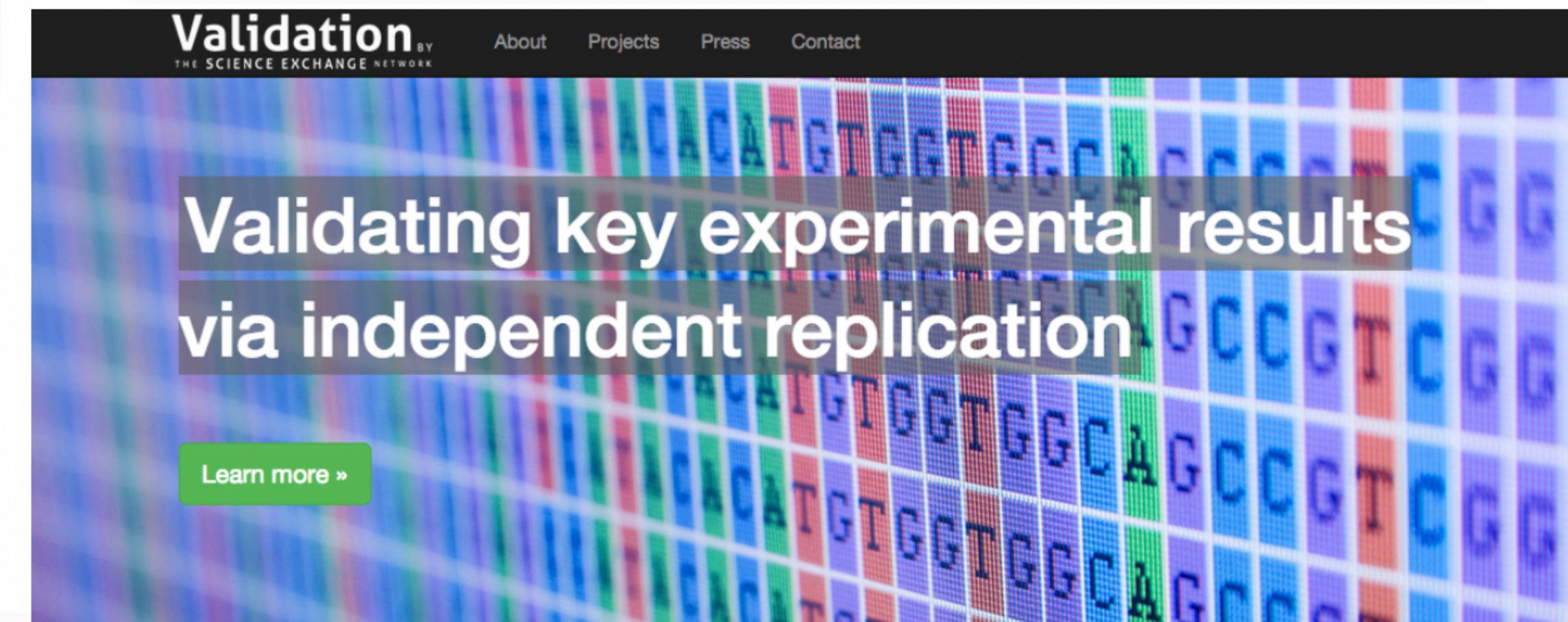
by Verne Kopytoff @vkopytoff JUNE 30, 2014, 10:58 AM EDT



PHYS.ORG

### Science is in a reproducibility crisis: How do we resolve it?

Sep 20, 2013 by Fiona Fidler & Ascelin Gordon, The Conversation



Validation BY THE SCIENCE EXCHANGE NETWORK

About Projects Press Contact

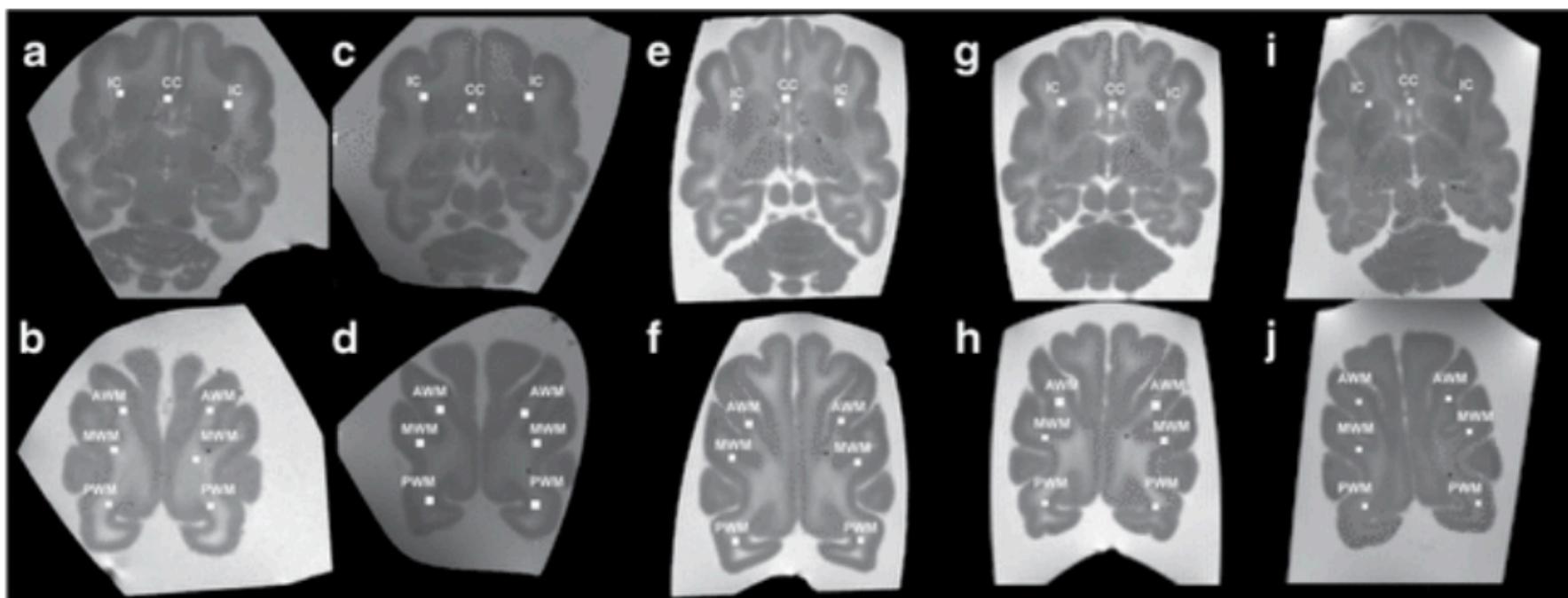
## Validating key experimental results via independent replication

Learn more »

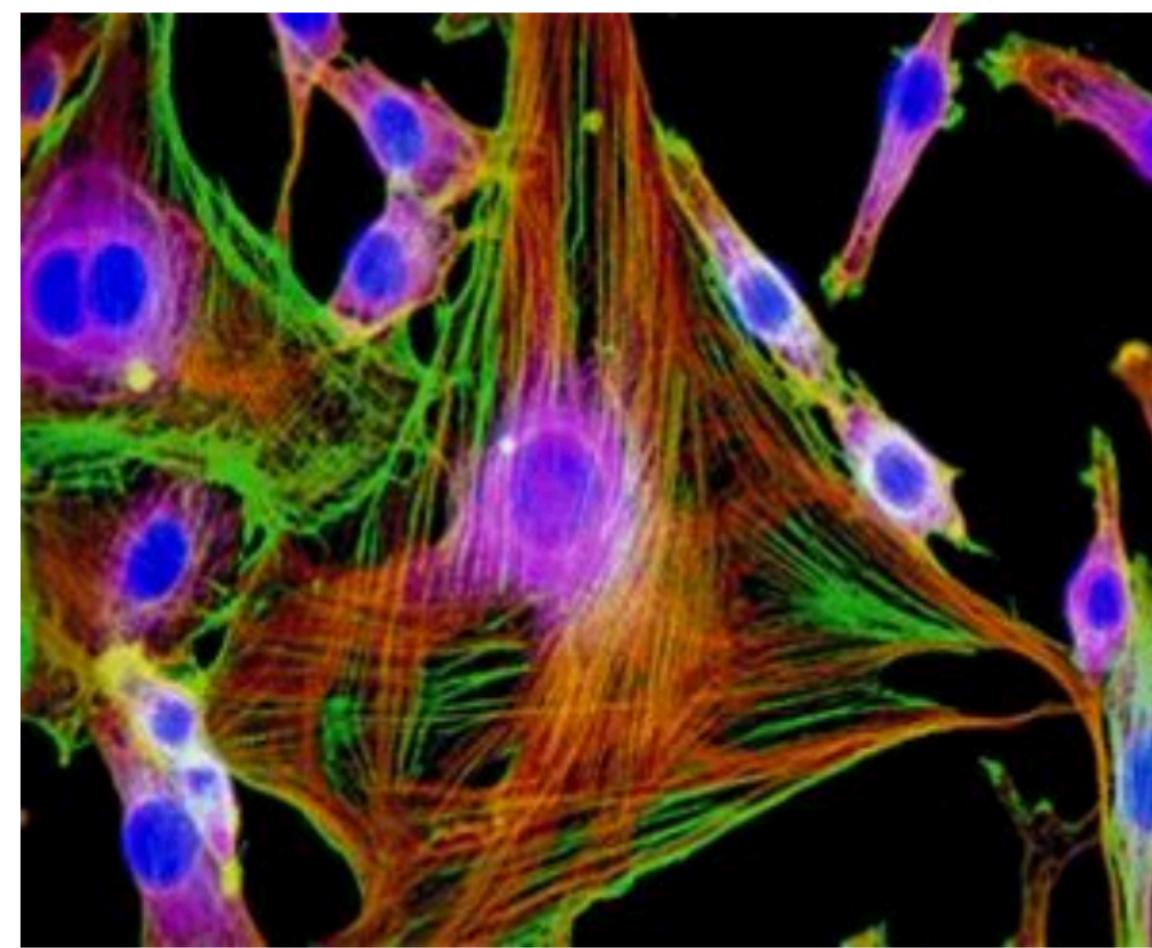
A background image of a DNA sequencing gel showing multiple lanes of sequence data, with the text overlaid on the right side.

Reproducibility initiatives will invite increased scrutiny into data cleaning methods.

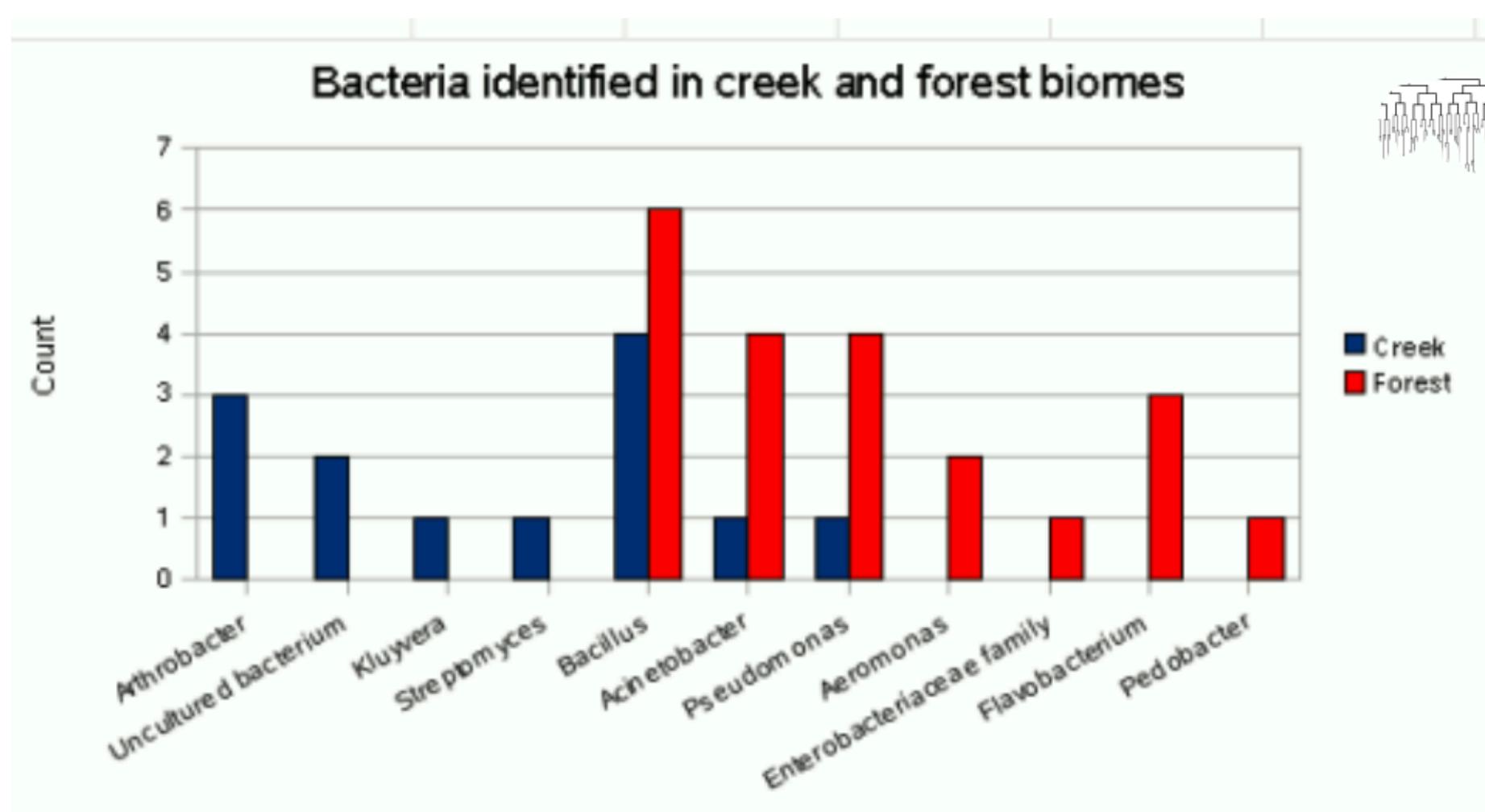
# Reproducible research



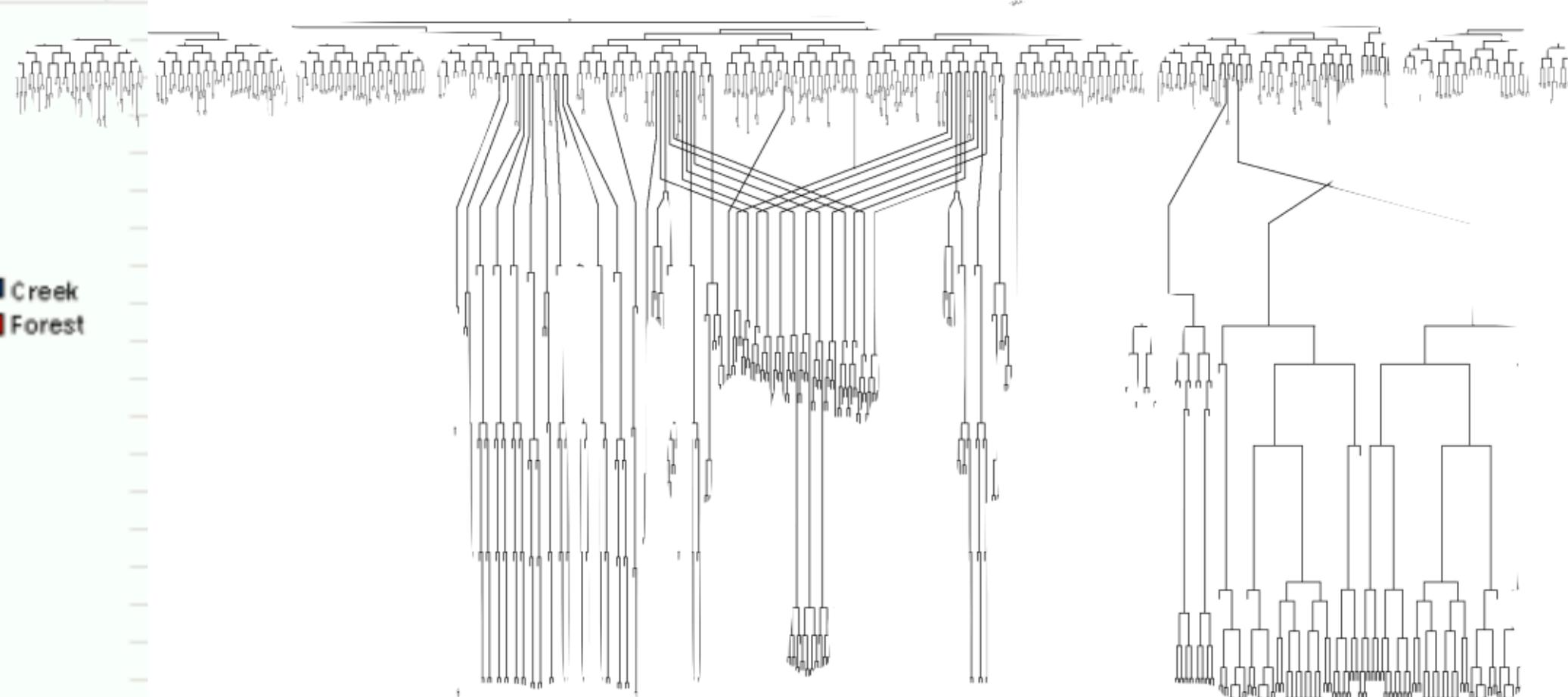
Reproducibility?



Sampling?



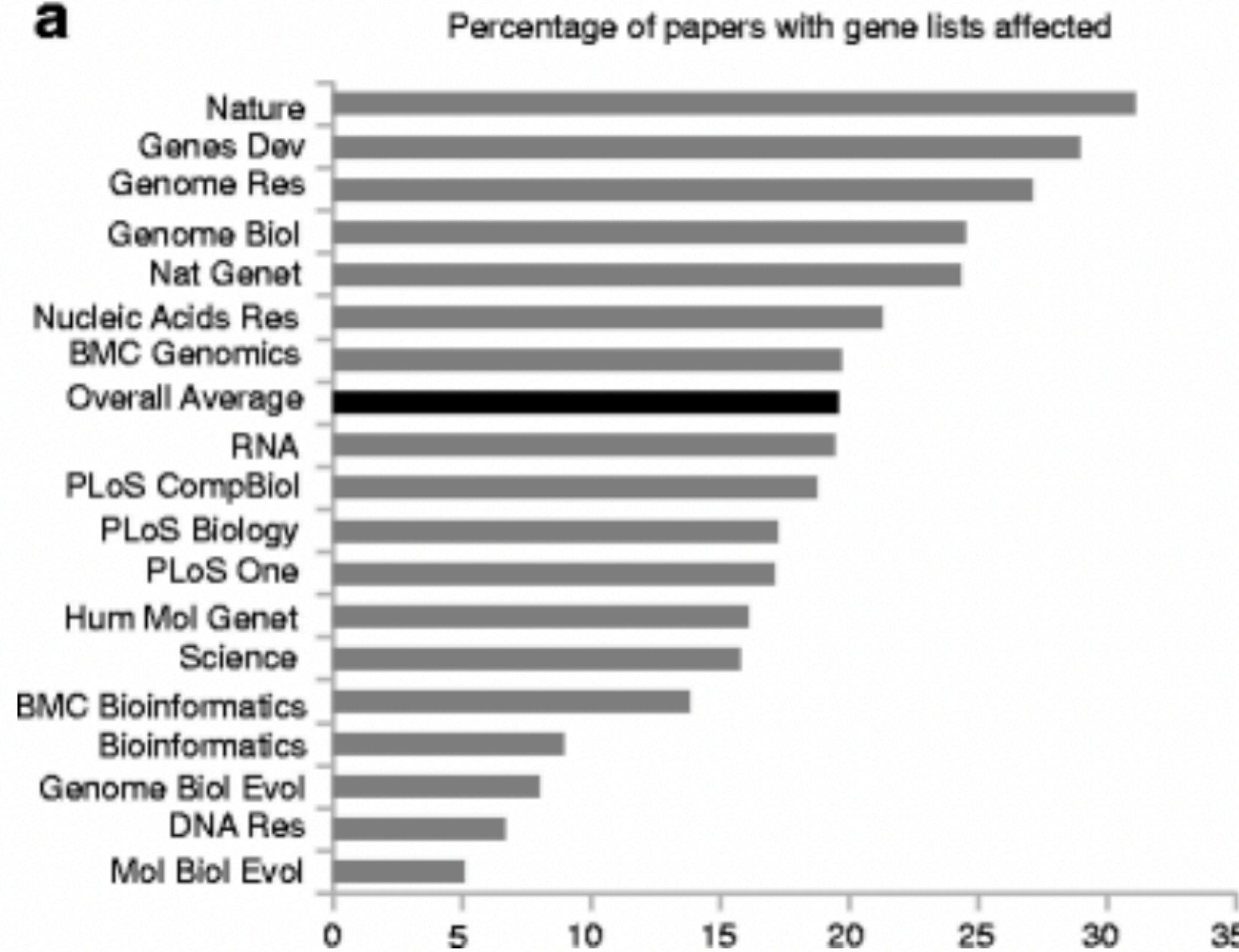
Error?



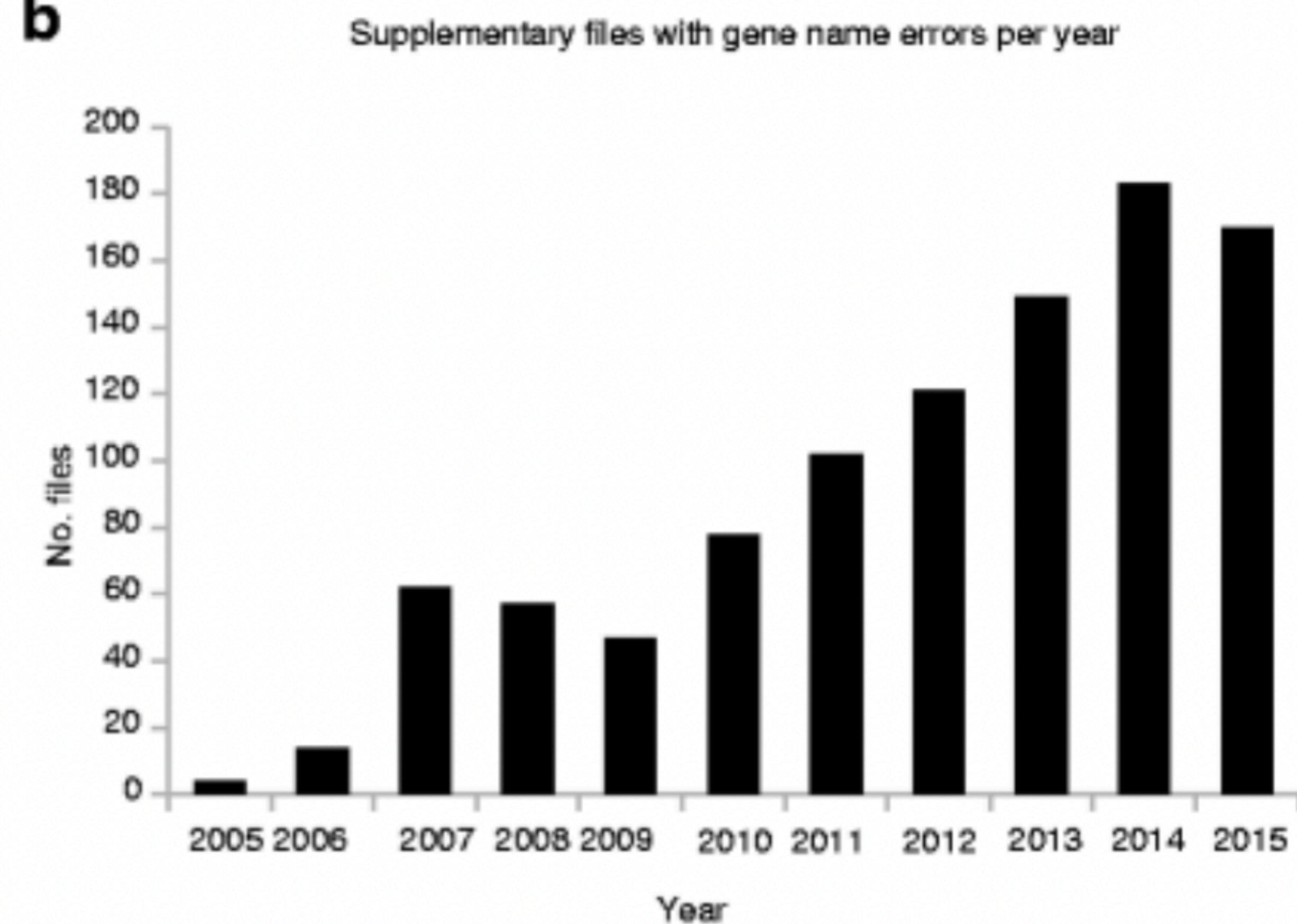
Automation?

# Reproducible research

a



b



# What I hope you learn from this class:

- Biology is becoming more and more **quantitative**, so data analysis, reproducible research, and proper statistics skills are more important than ever
- Excel is great for some things, but data analysis should be **scripted** (in R or another language) in a way that any person could take your raw data and reproduce the figures in your manuscript
- Statistics can be slimy! (1) you must be **cautious and critical** when reviewing others' statistics and (2) you must be **transparent and honest** when providing your own analysis

# Science, variability, and statistics

***Response of sheep to anthrax***

Response	Vaccinated	Not vaccinated
Died of anthrax	0	24
Survived	24	0
Total	24 (100%)	0 (0%)

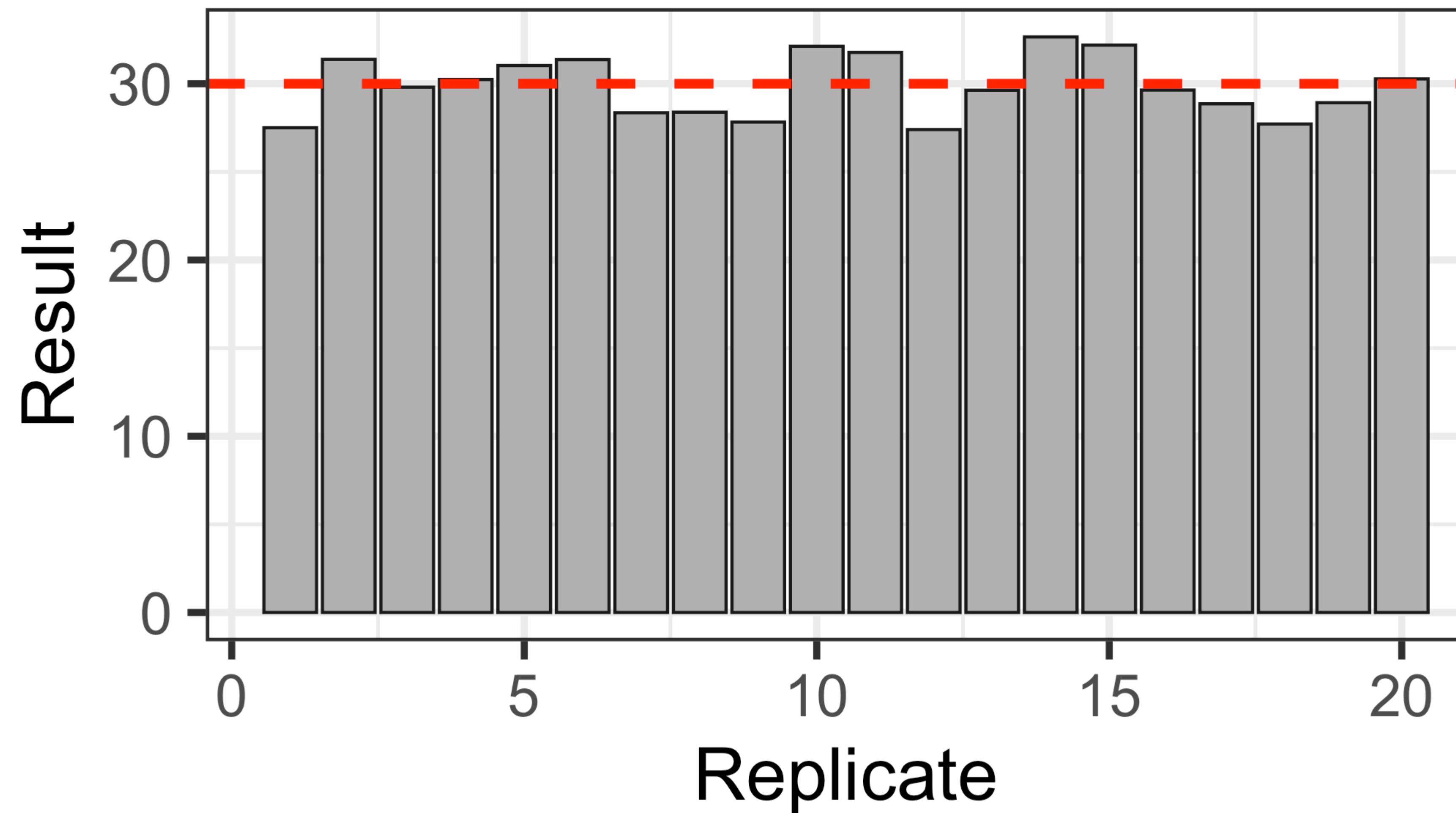
***Incidence of liver tumors in mice***

Response	<i>E. coli</i>	Germ free
Liver tumors	8	19
No tumors	5	30
Total	13 (62%)	49 (39%)

✓ ***Significant!***

⚠ ? ***Significant?***

# Science, variability, and statistics



# Science, variability, and statistics

***Response of sheep to anthrax***

Response	Vaccinated	Not vaccinated
Died of anthrax	0	3 <del>24</del>
Survived	3 <del>24</del>	0
Total	3 <del>24</del> (100%)	0 (0%)

***Incidence of liver tumors in mice***

Response	E. coli	Germ free
Liver tumors	8	19
No tumors	5	30
Total	13 (62%)	49 (39%)



**Significant?**



**Significant?**

# Observational studies v. experiments

**Collecting data from subjects  
as an observer, not  
manipulating conditions**

Survey about smoking habits  
and health

Performing a GWA/disease risk  
analysis on individuals with  
different genotypes

**Collecting data in a  
controlled environment  
where researchers impose  
the conditions**

Clinical trial for a new drug  
treatment

Knocking out a gene and then  
testing growth/survival with  
and without the gene

# Controls: an essential part of any study

**Collecting data from subjects  
as an observer, not  
manipulating conditions**

(Smokers & non-smokers)

Survey about smoking habits  
and health

(Disease & non-disease)

Performing a GWA/disease risk  
analysis on individuals with  
different genotypes

**Collecting data in a  
controlled environment  
where researchers impose  
the conditions**

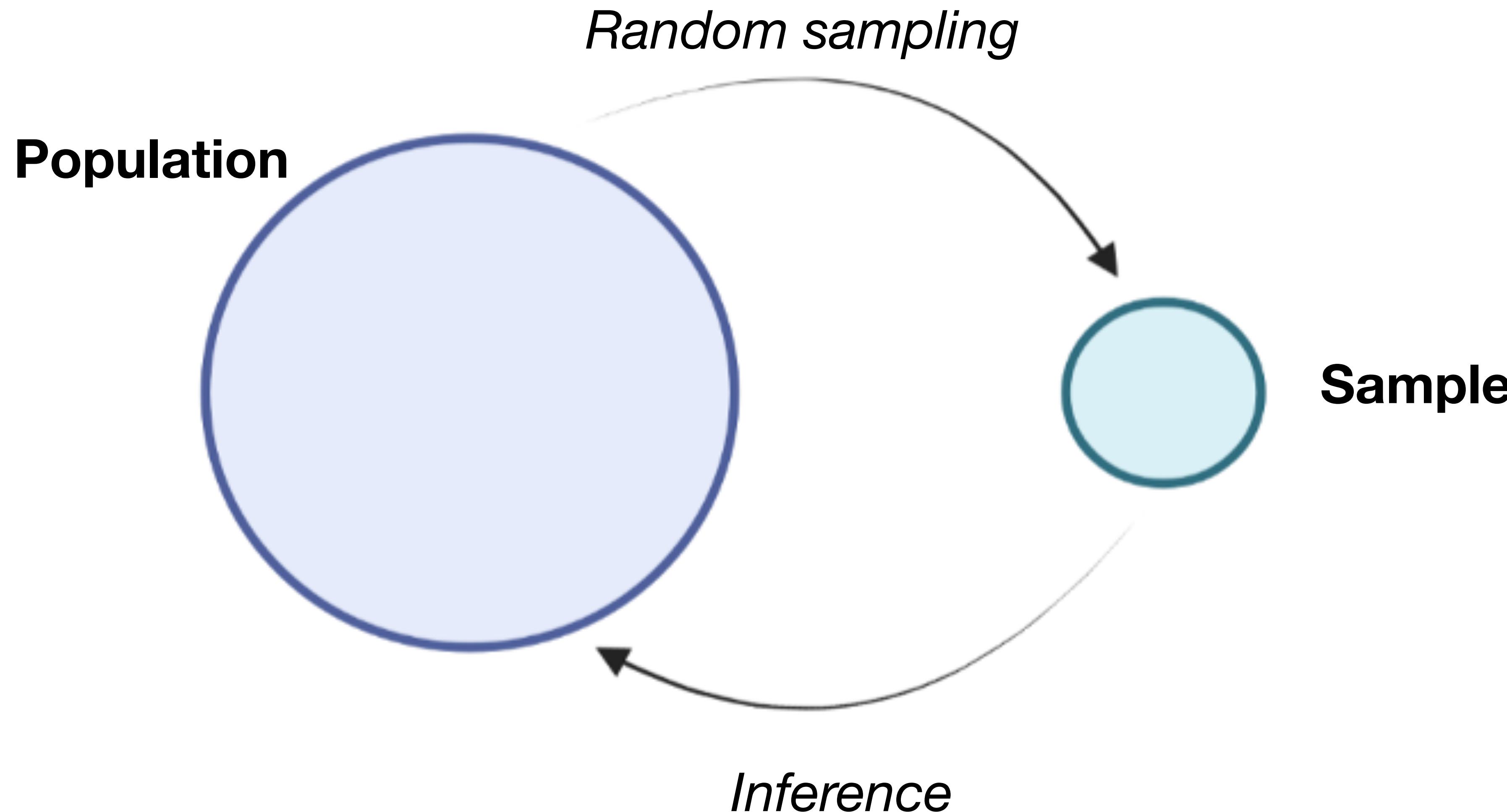
(Drug & placebo)

Clinical trial for a new drug  
treatment

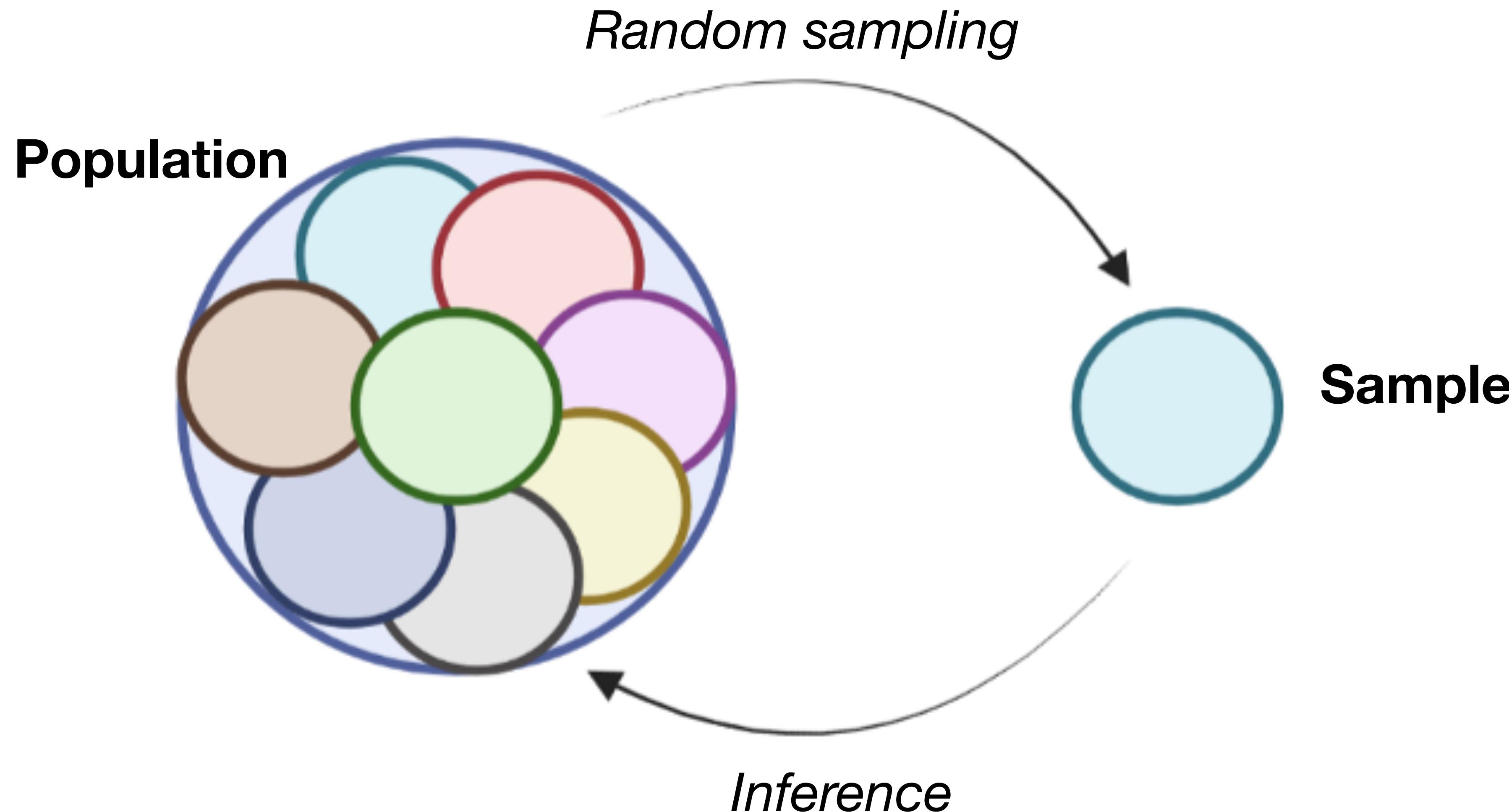
(Knockout & wild-type)

Knocking out a gene and then  
testing growth/survival with  
and without the gene

# Populations and sampling



# Populations and sampling



# Simple random sample

- Every member of the population has the **same** chance of being included in the sample
- Members of the sample are chosen **independently** of each other
  - *Not dependent on which other members are chosen*
- How we gather our data has **tremendous** implications on our choice of analysis methods and validity of our studies
  - ***There is no replacement for good/clean data!!!***

# Types of variables

## Categorical

Blood type (A, B, AB, O)

Fish sex (male, female)

Shape of pea (smooth, wrinkled)

Success in trial (Alive, dead)

## Numeric

Human height

Blood cholesterol of patient

Number of bacterial colonies

Length of DNA segment

# Types of variables

# Discrete

# Categorical

# Blood type (A, B, AB, O)

# Fish sex (male, female)

# Shape of pea (smooth, wrinkled)

# Success in trial (Alive, dead)

# ***Continuous***

# Numeric

# *Human height*

# *Blood cholesterol of patient*

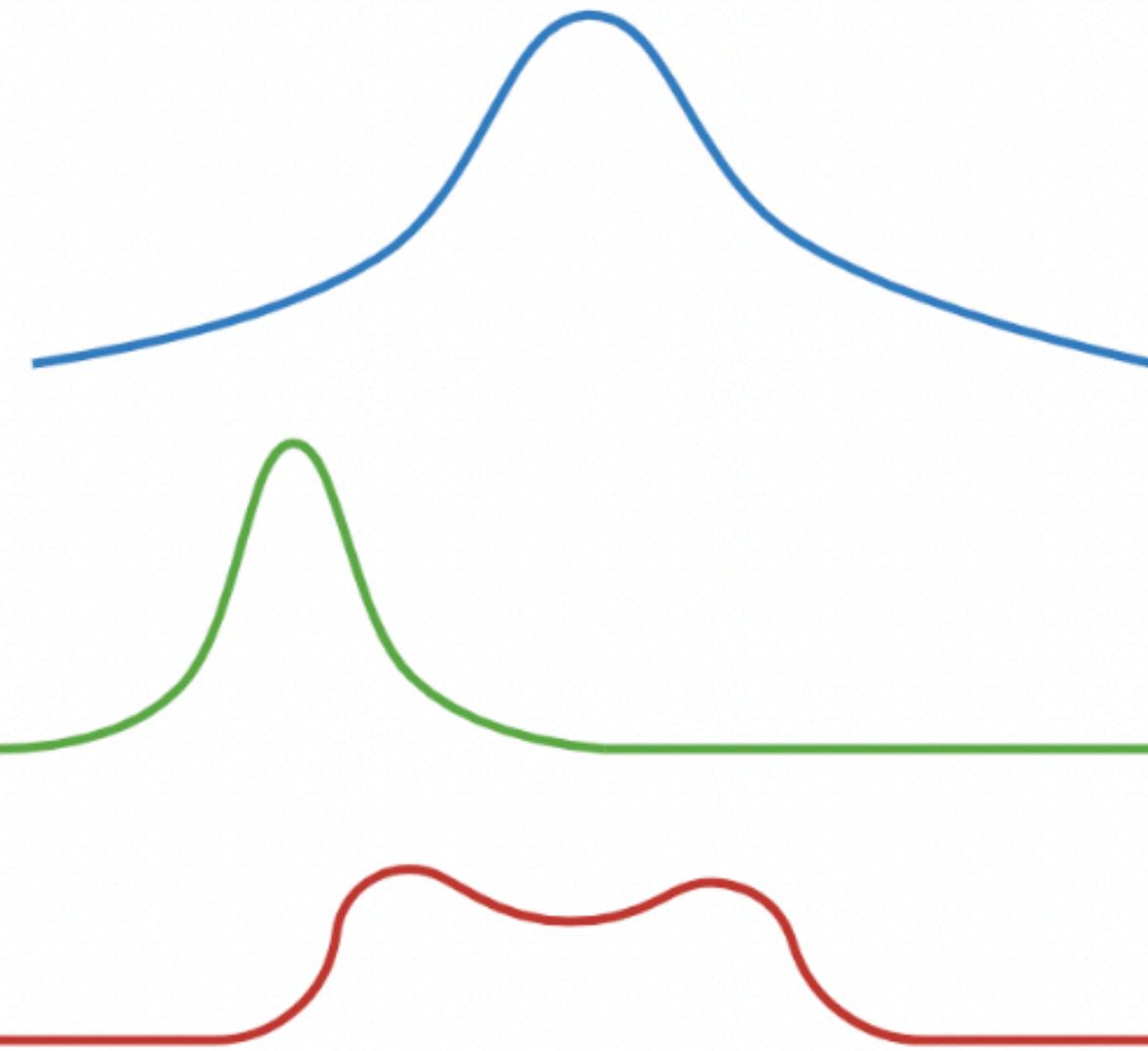
# Number of bacterial colonies

# Length of DNA segment

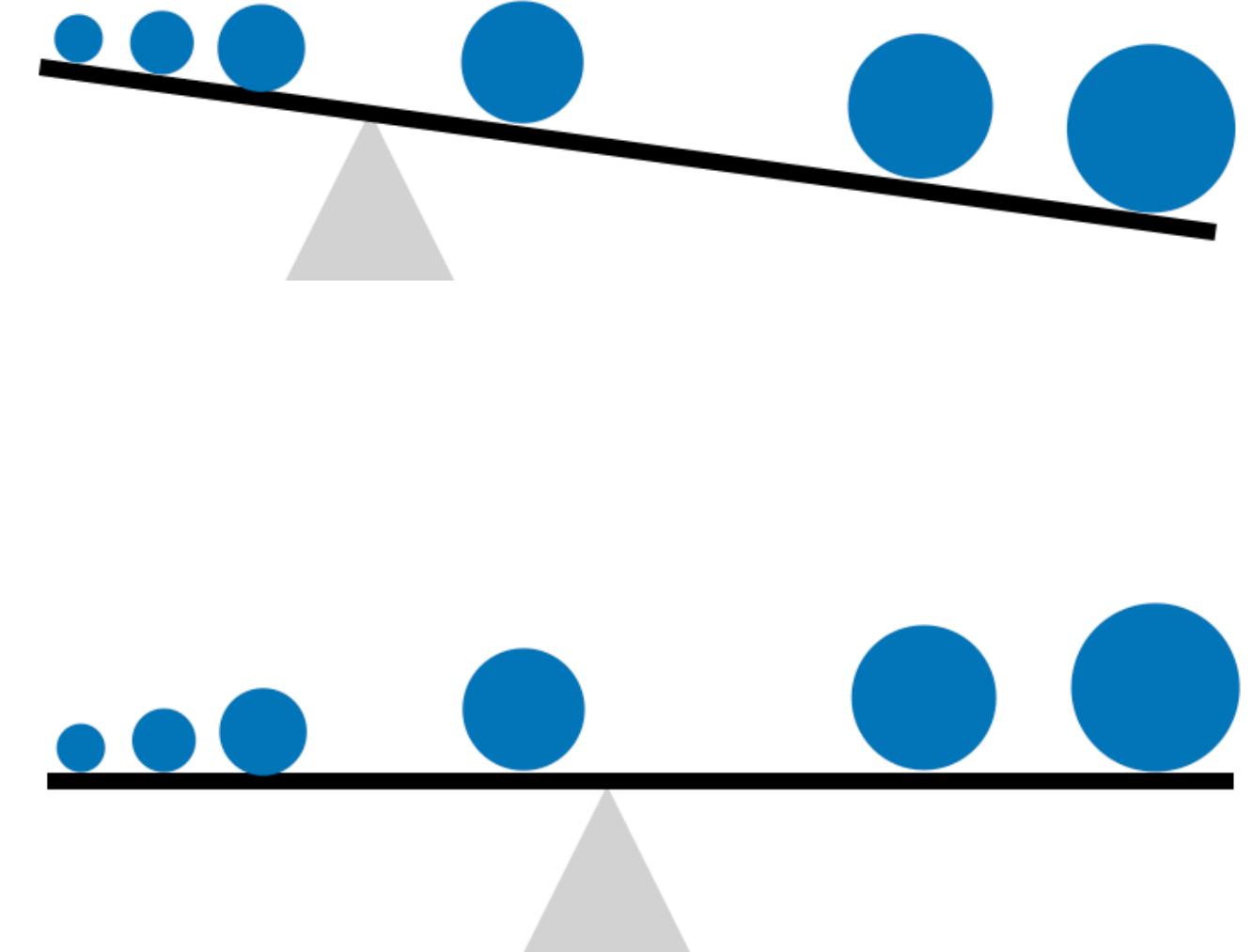
# Discrete

# Data exploration and summarization

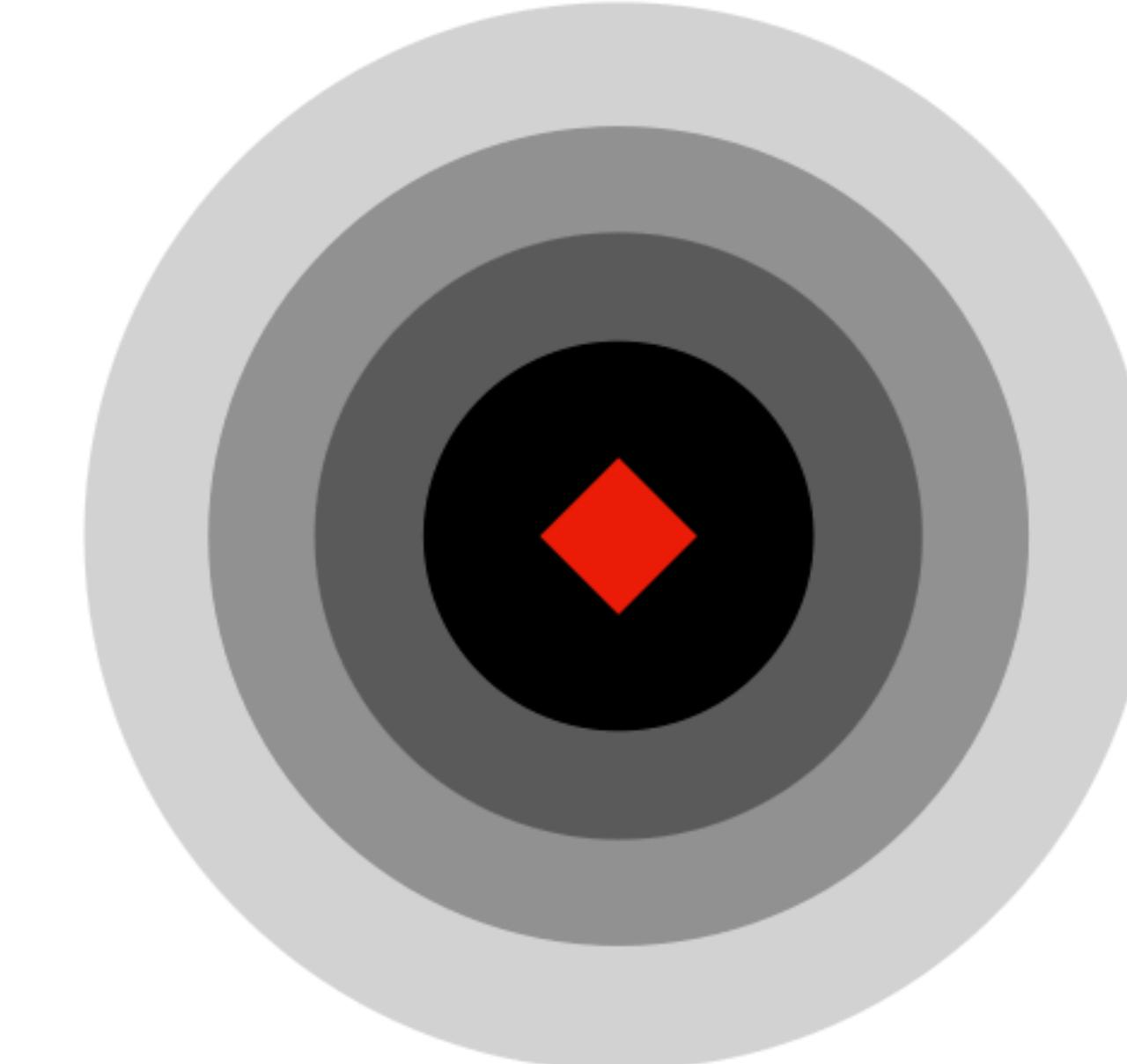
Shape



Center

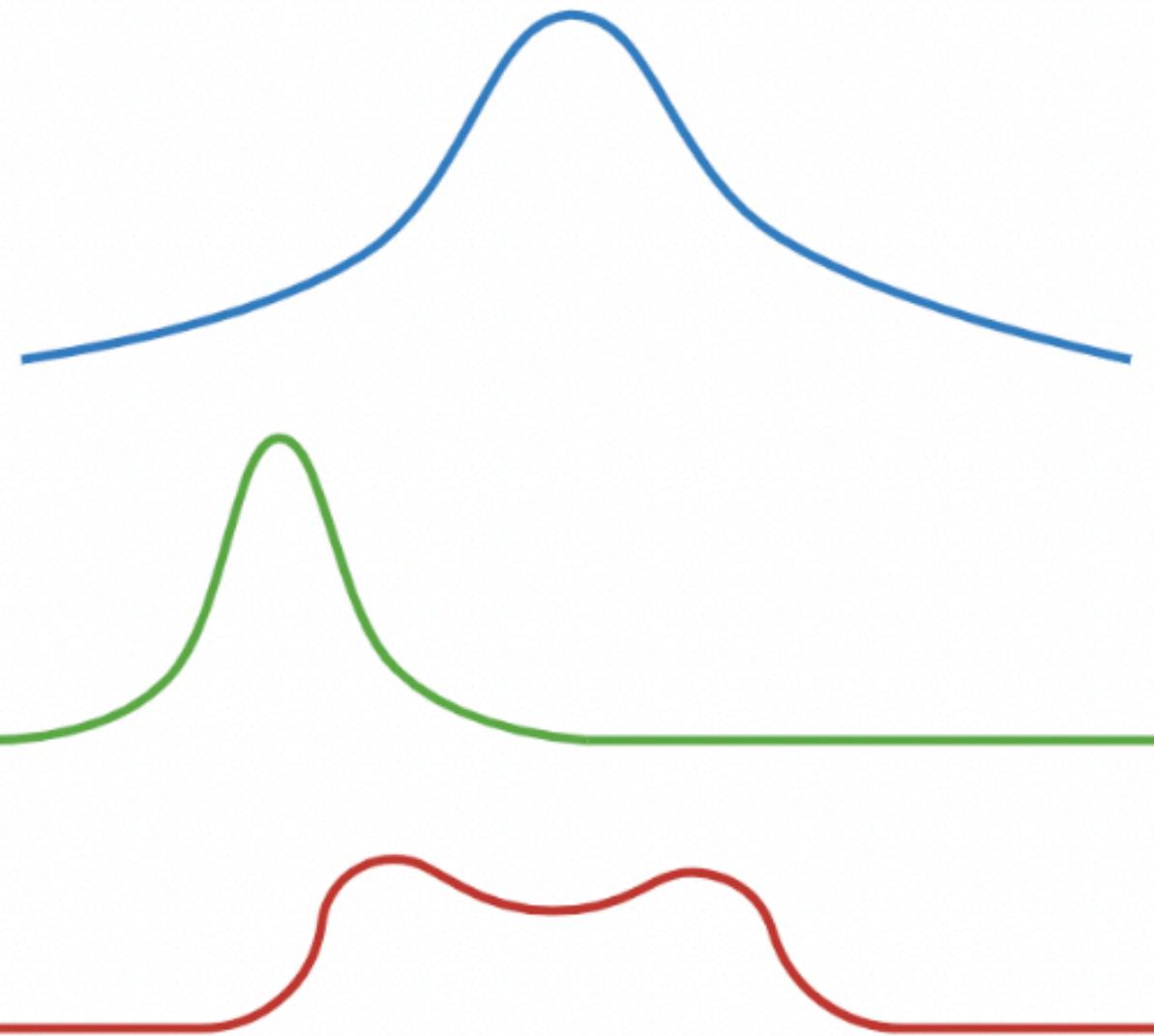


Spread

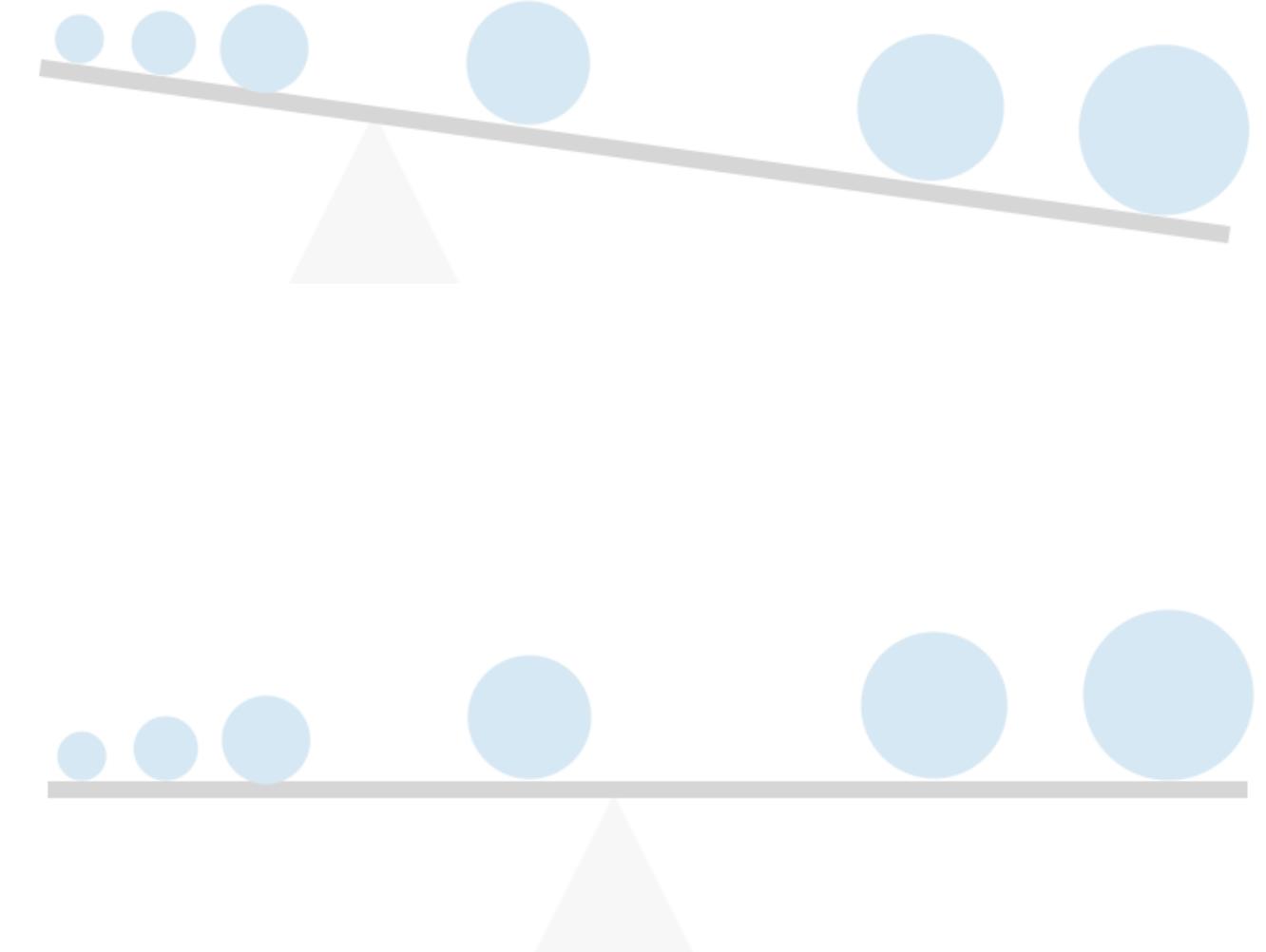


# Data exploration and summarization

Shape



Center



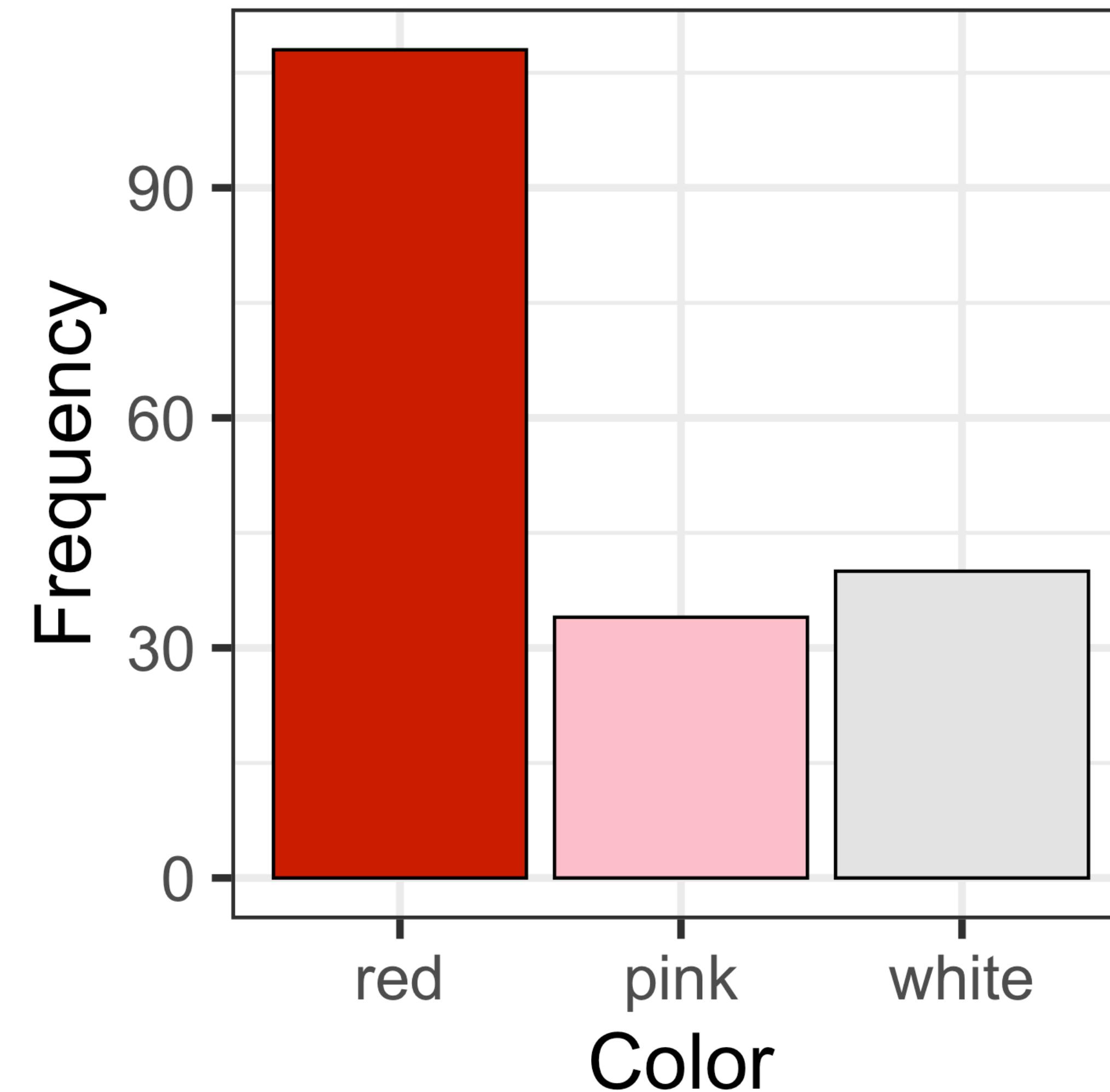
Spread



# Visualizing frequency distributions

Color	Frequency	Relative frequency	Percent frequency
Red	108	0.59	59.34
Pink	34	0.19	18.68
White	40	0.22	21.98
Total	182	1	100

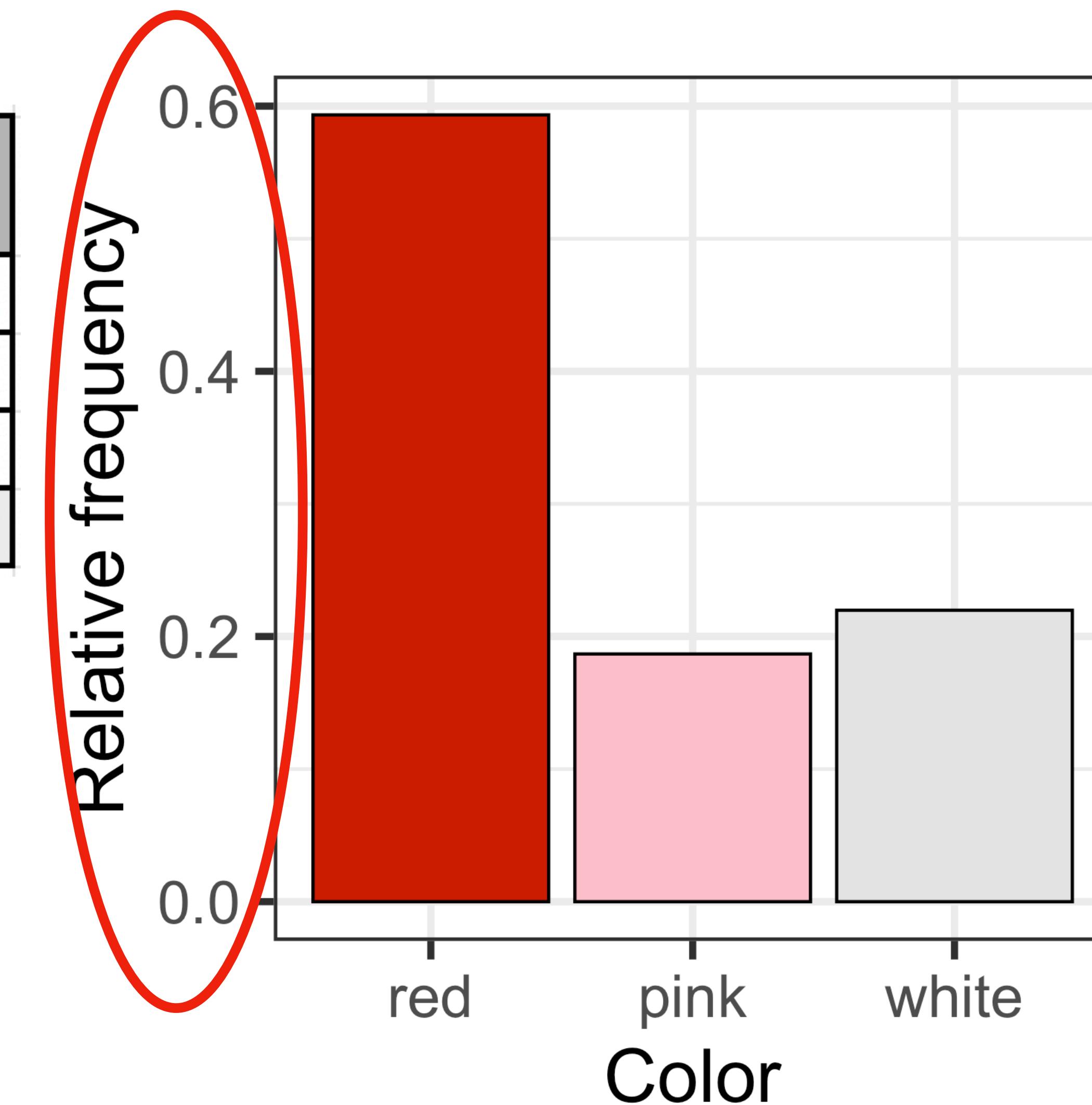
- Categorical data can be represented with a **bar chart**



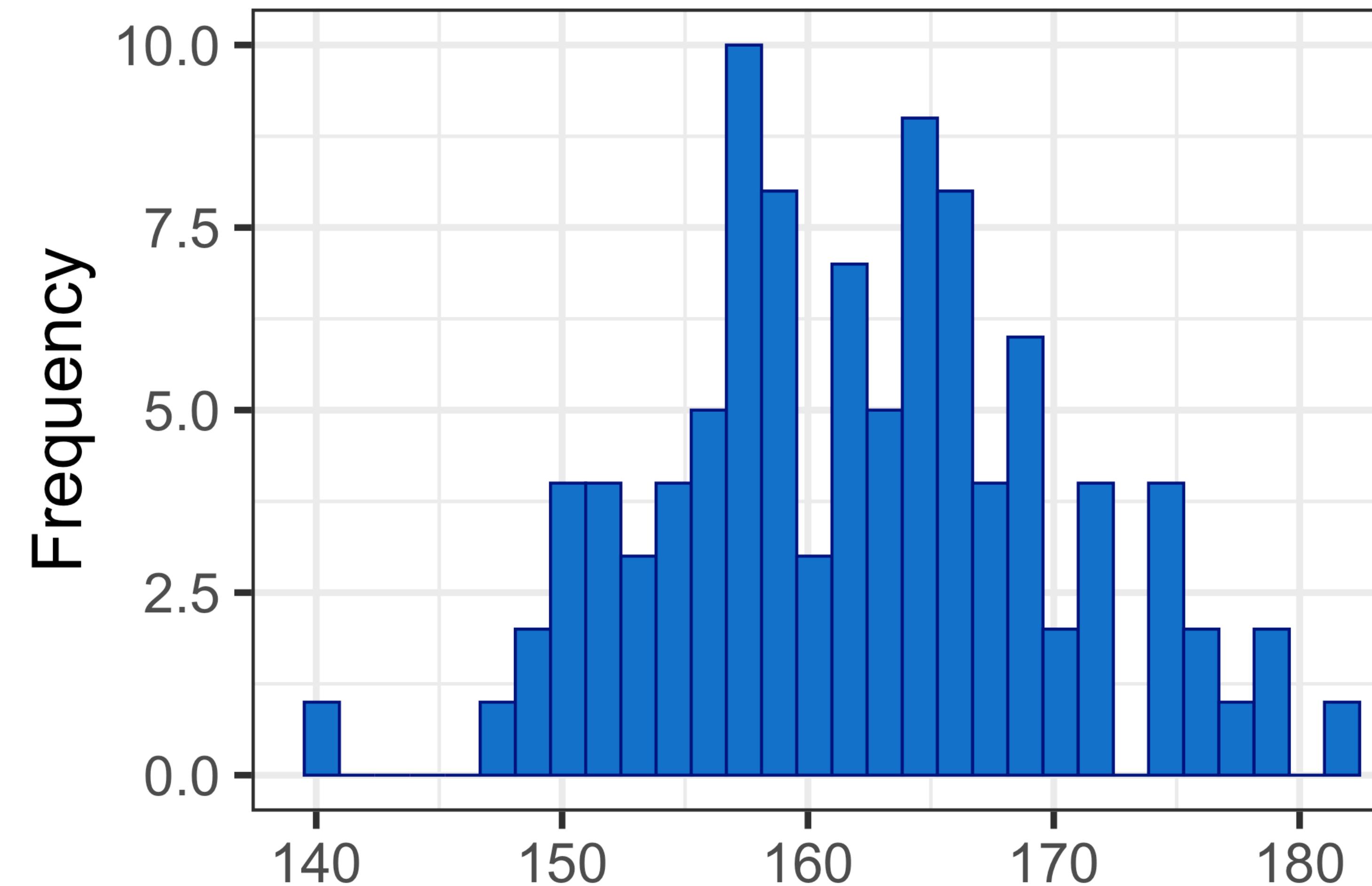
# Visualizing frequency distributions

Color	Frequency	Relative frequency	Percent frequency
Red	108	0.59	59.34
Pink	34	0.19	18.68
White	40	0.22	21.98
Total	182	1	100

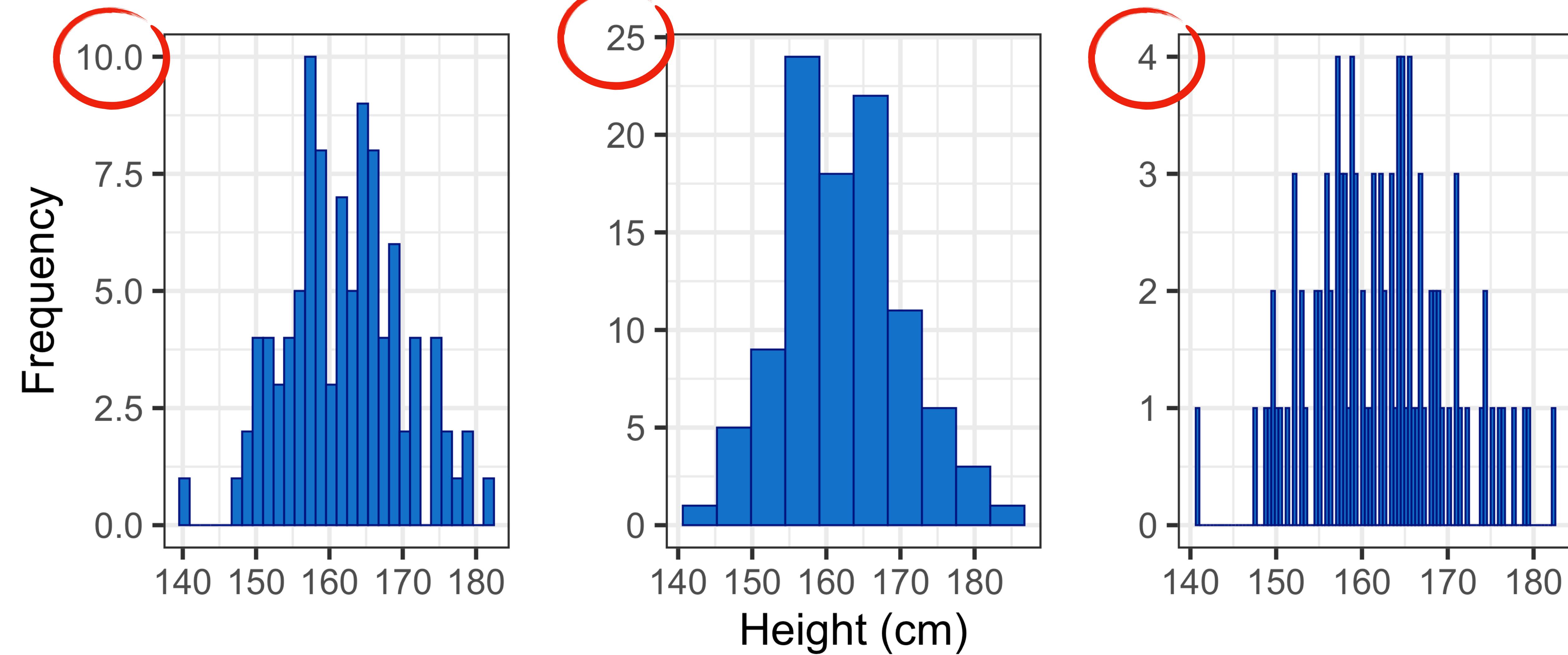
- Categorical data can be represented with a **bar chart**
- **Relative frequency** can be useful to compare datasets of different sizes



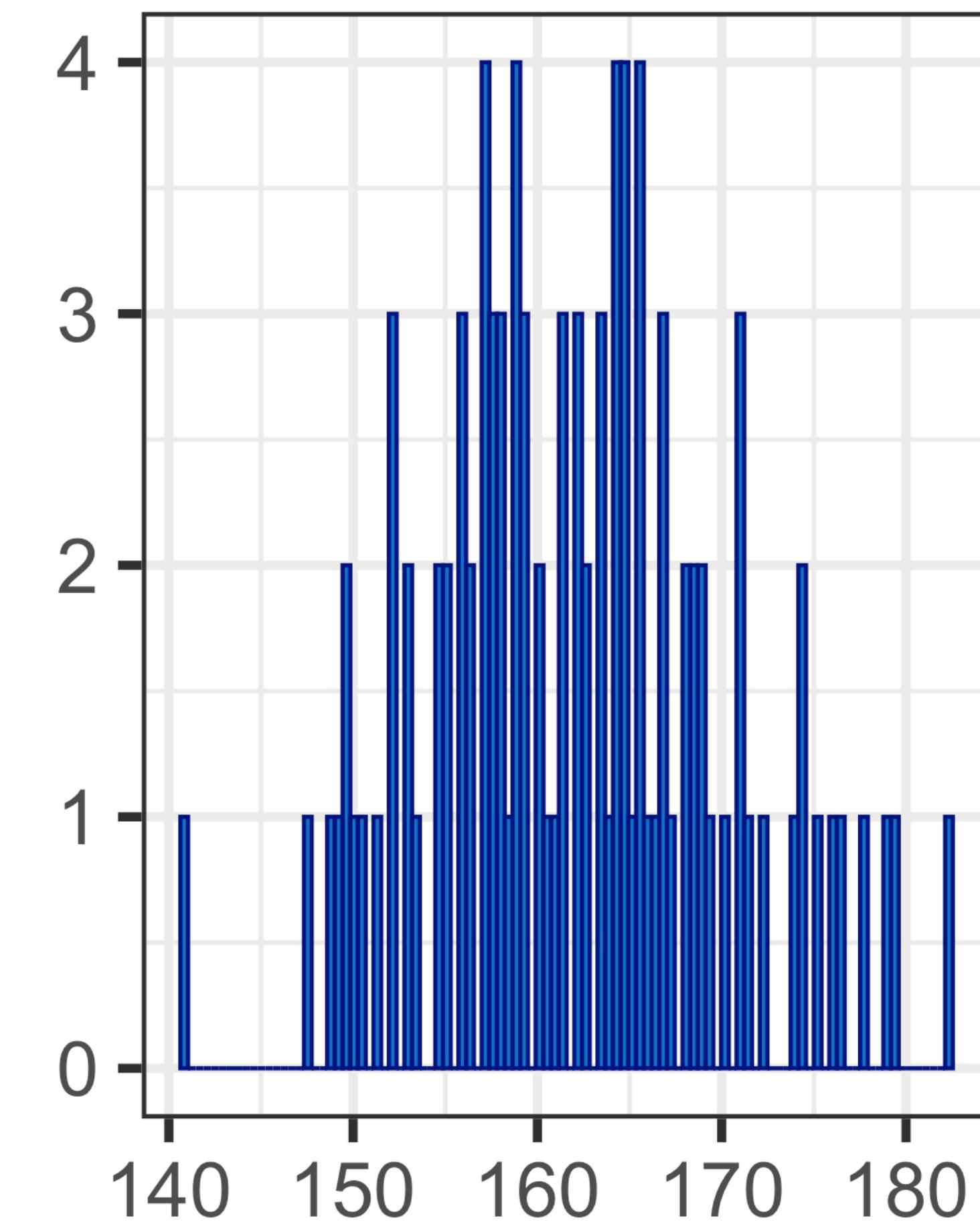
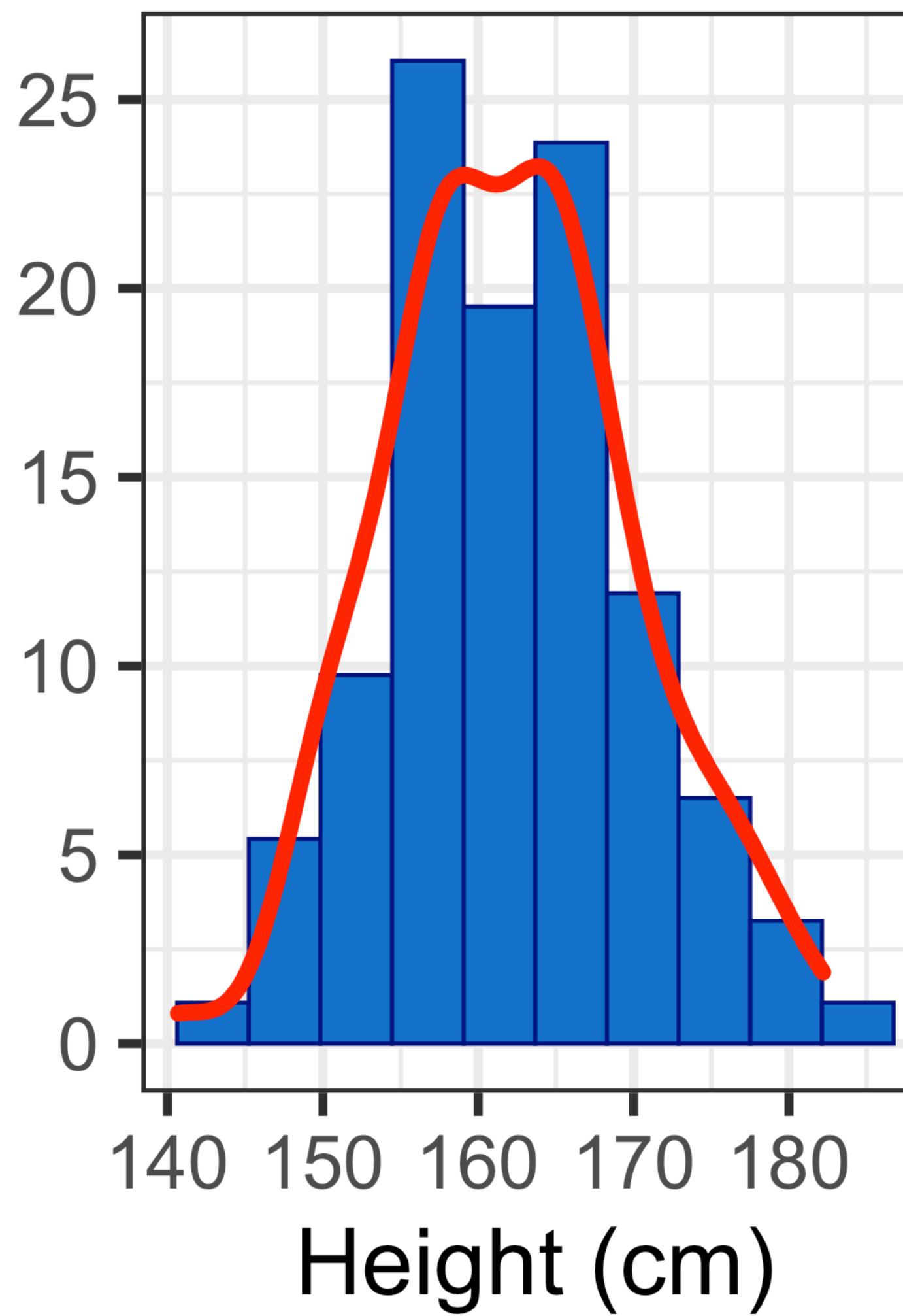
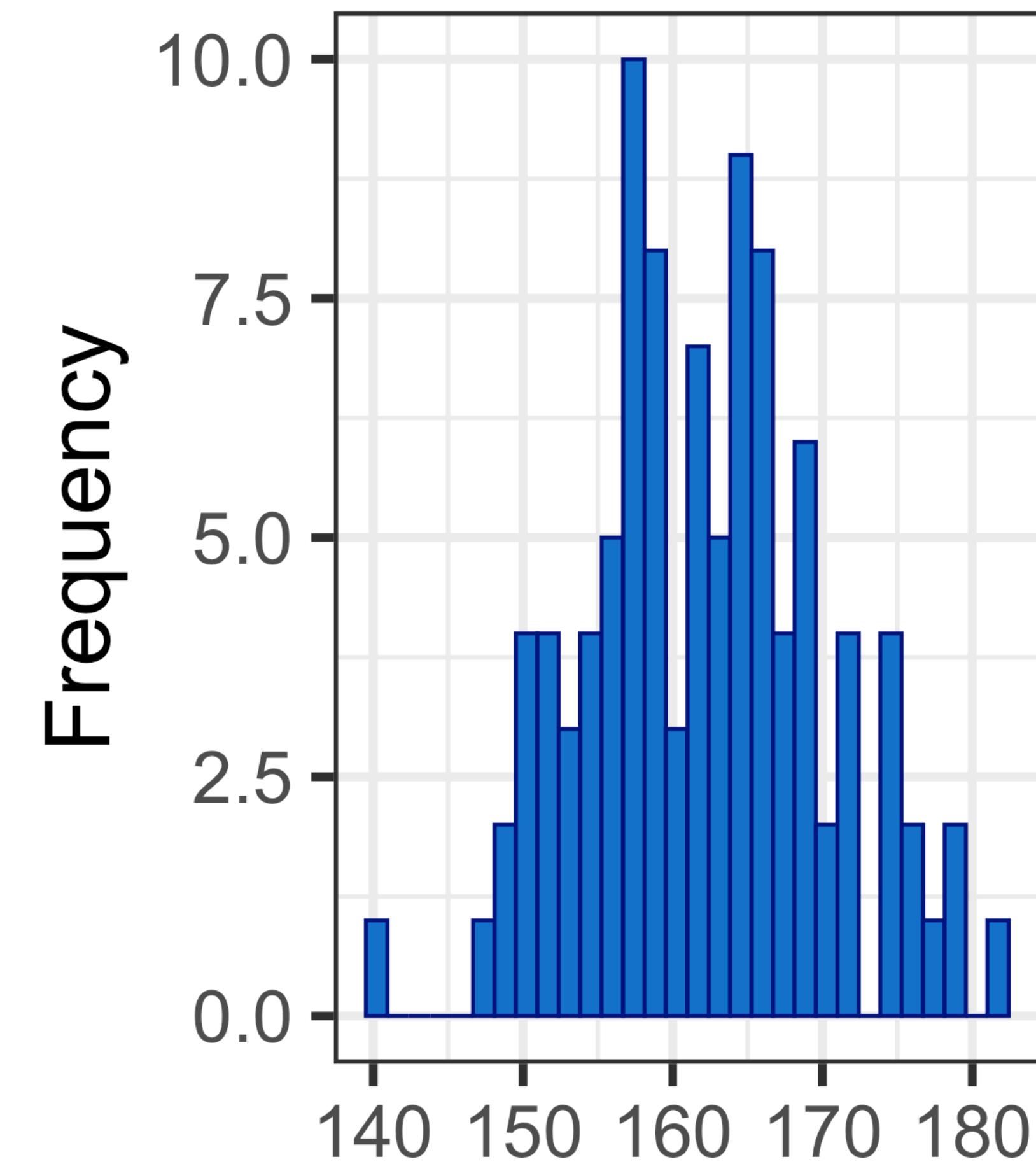
# Histograms show shape of numerical data



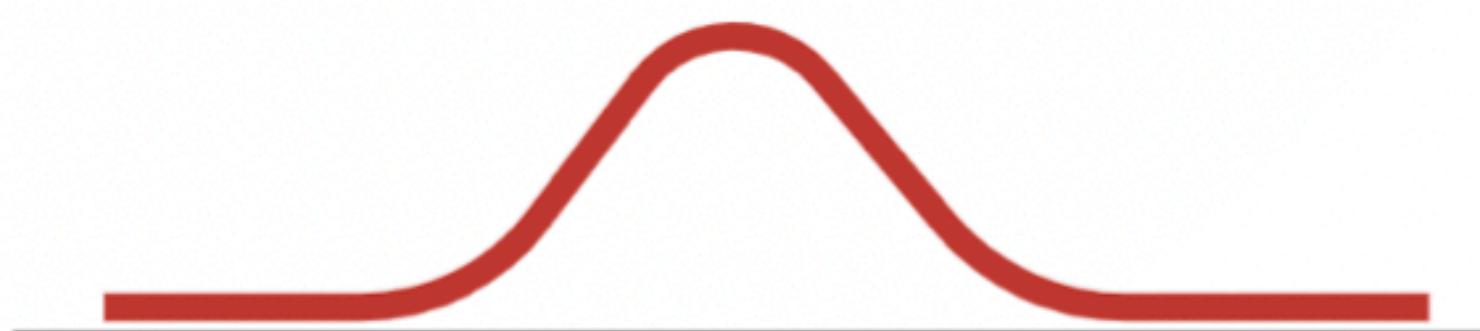
# Histograms show shape of numerical data



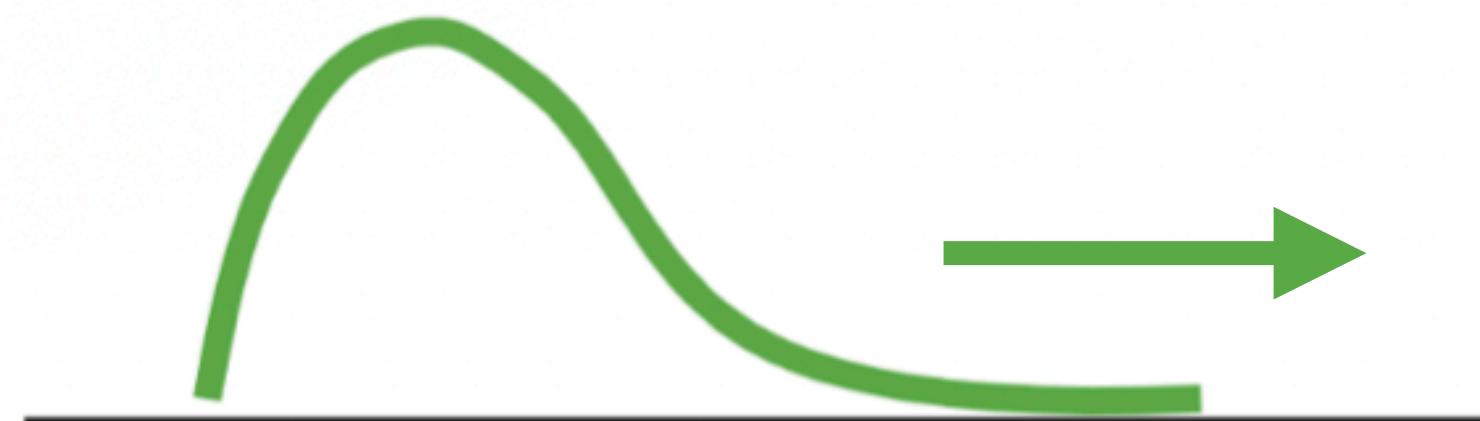
# Histograms show shape of numerical data



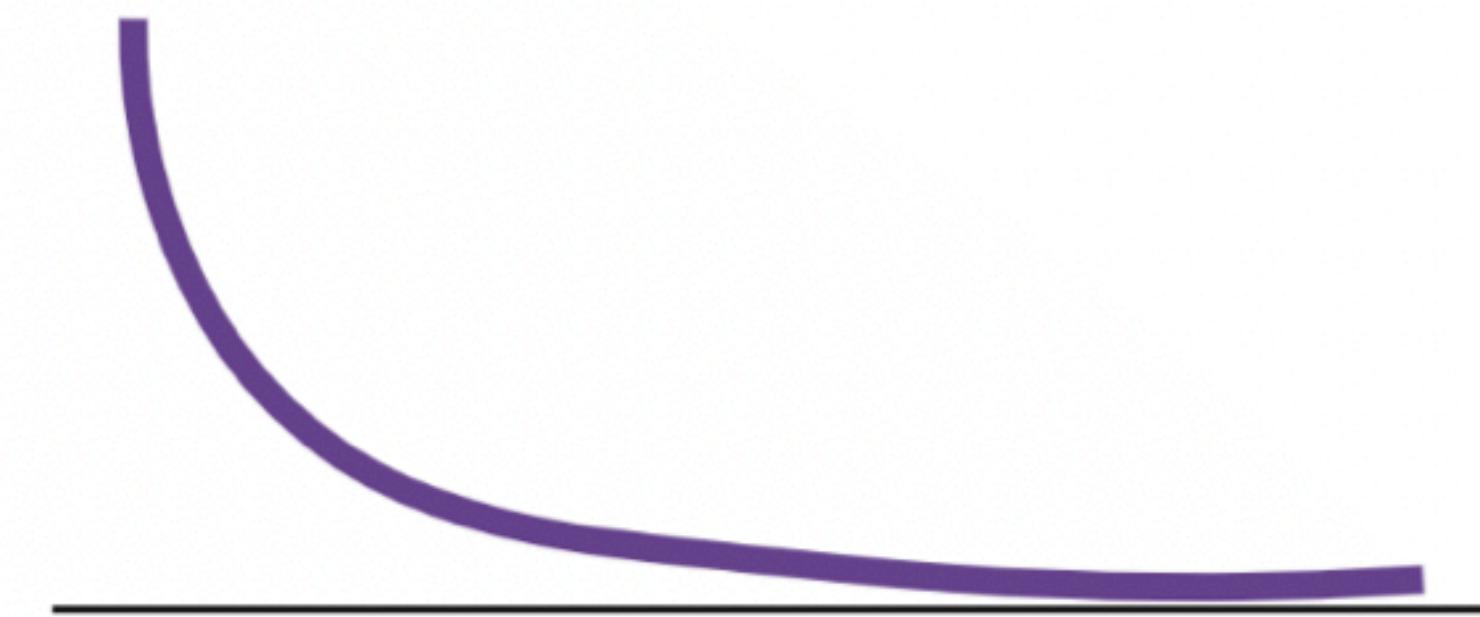
# Biological frequency distribution shapes



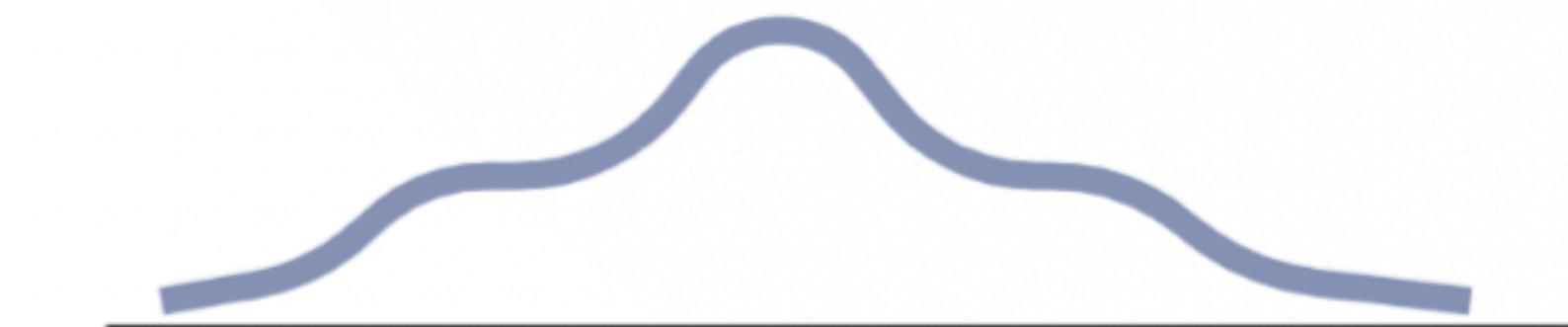
Symmetric, bell-shaped



Skewed to the right



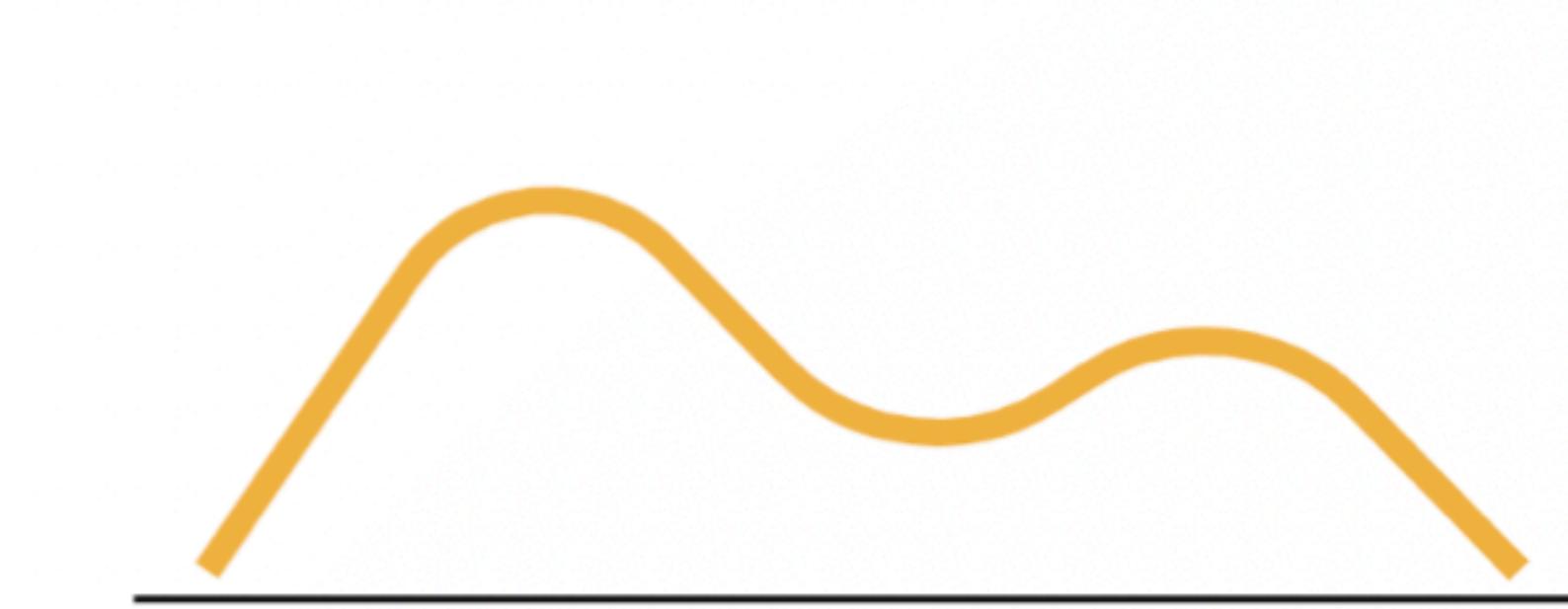
Exponential



Symmetric, not bell-shaped



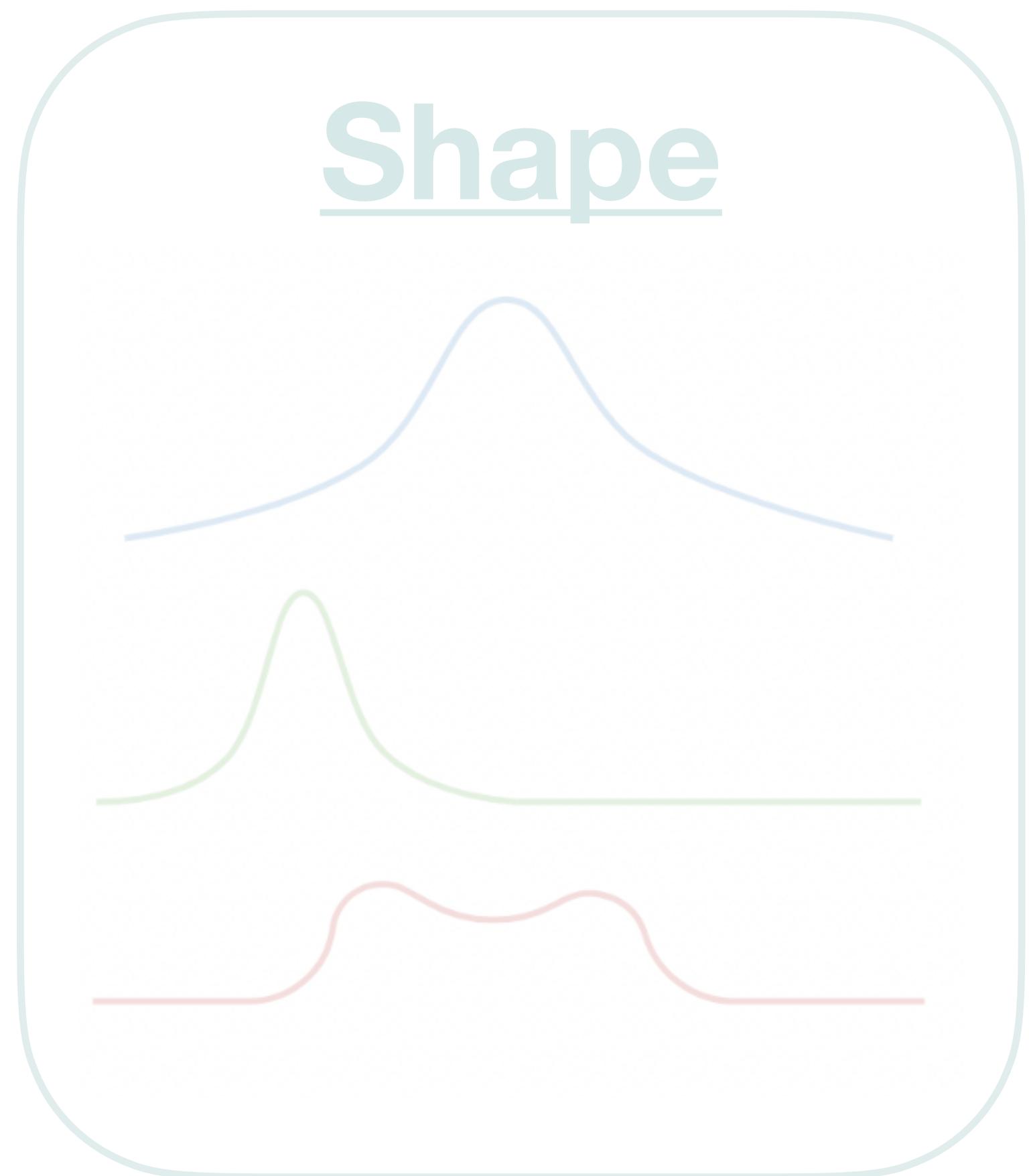
Skewed to the left



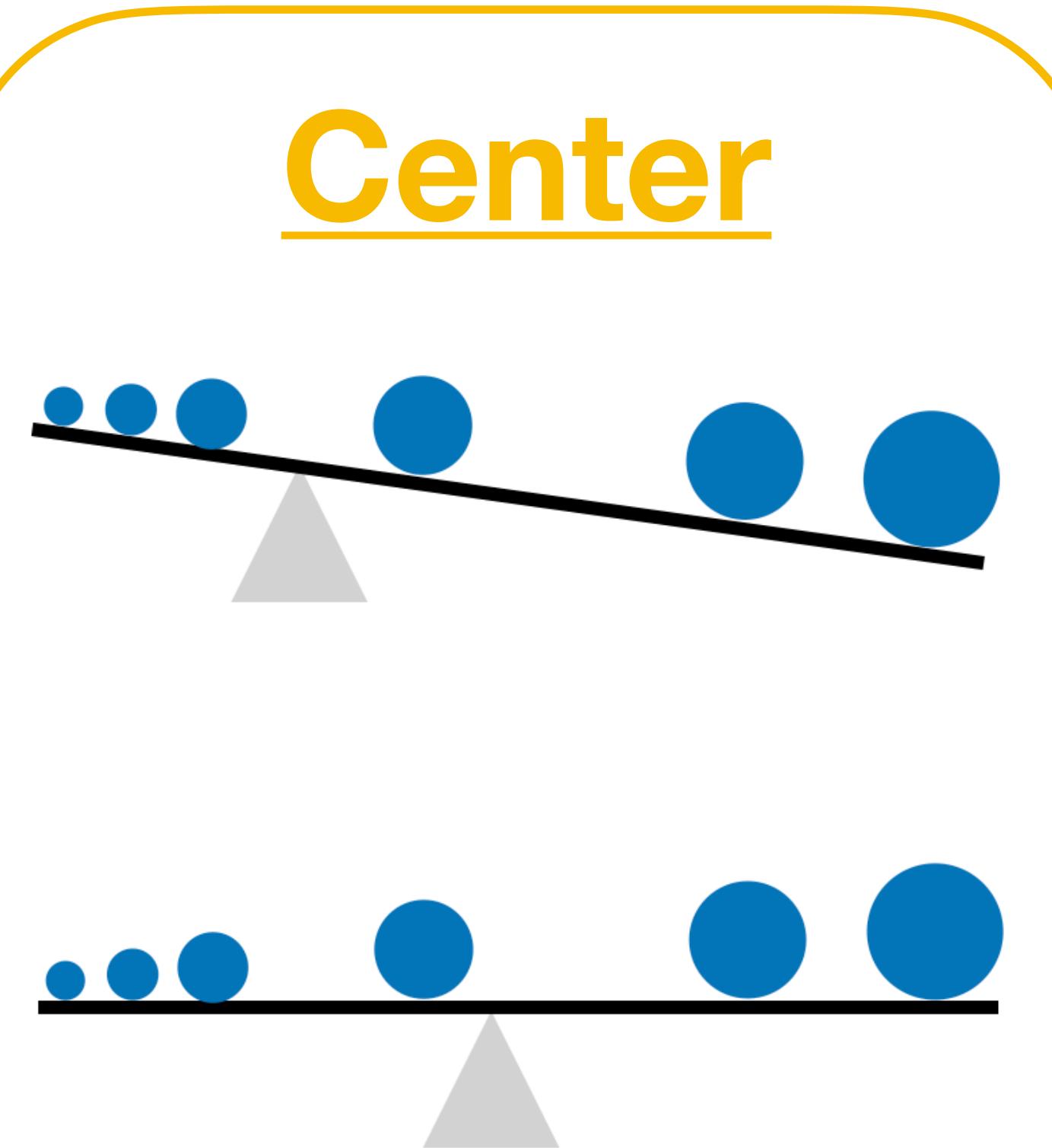
Bimodal

# Data exploration and summarization

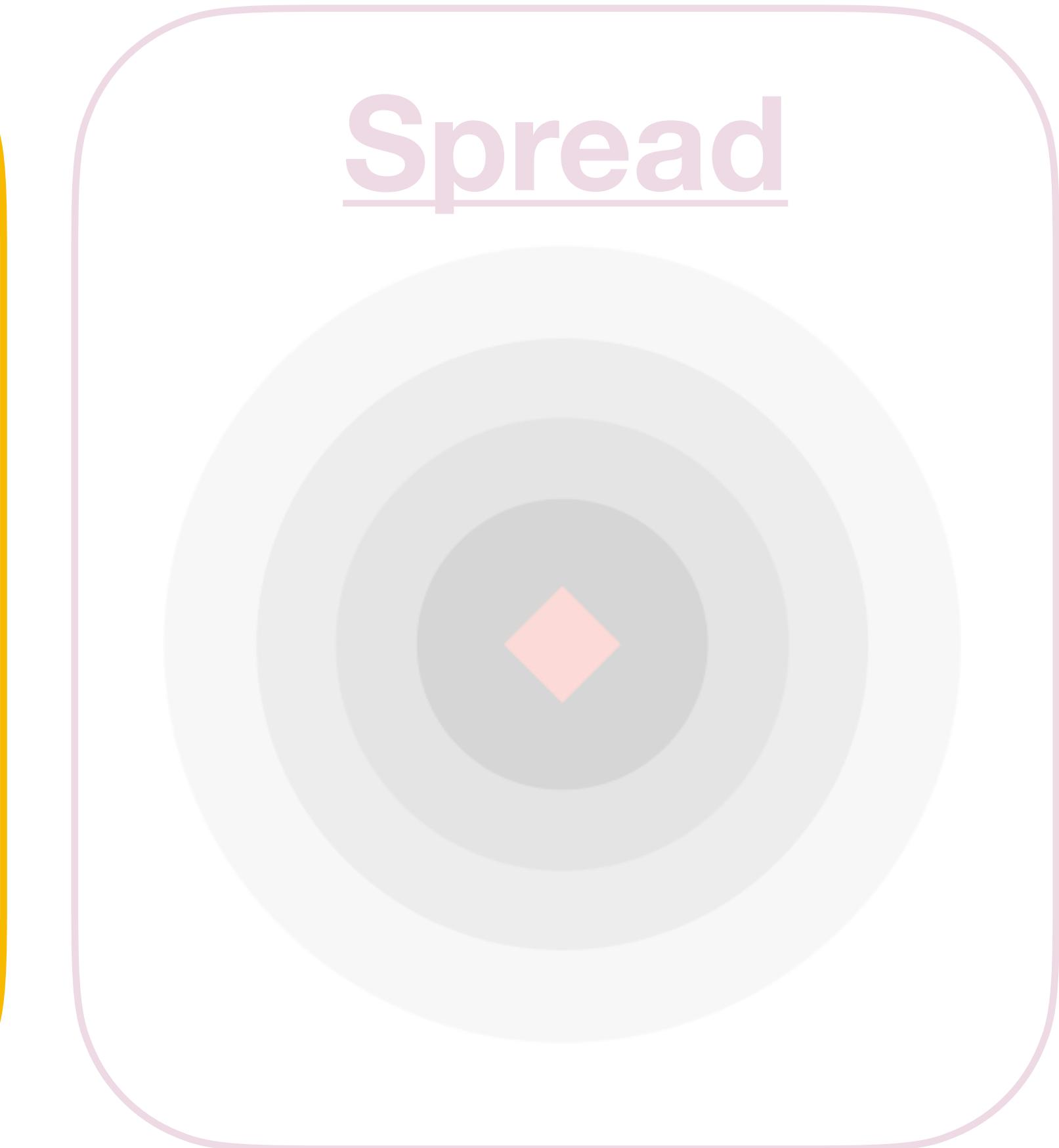
Shape



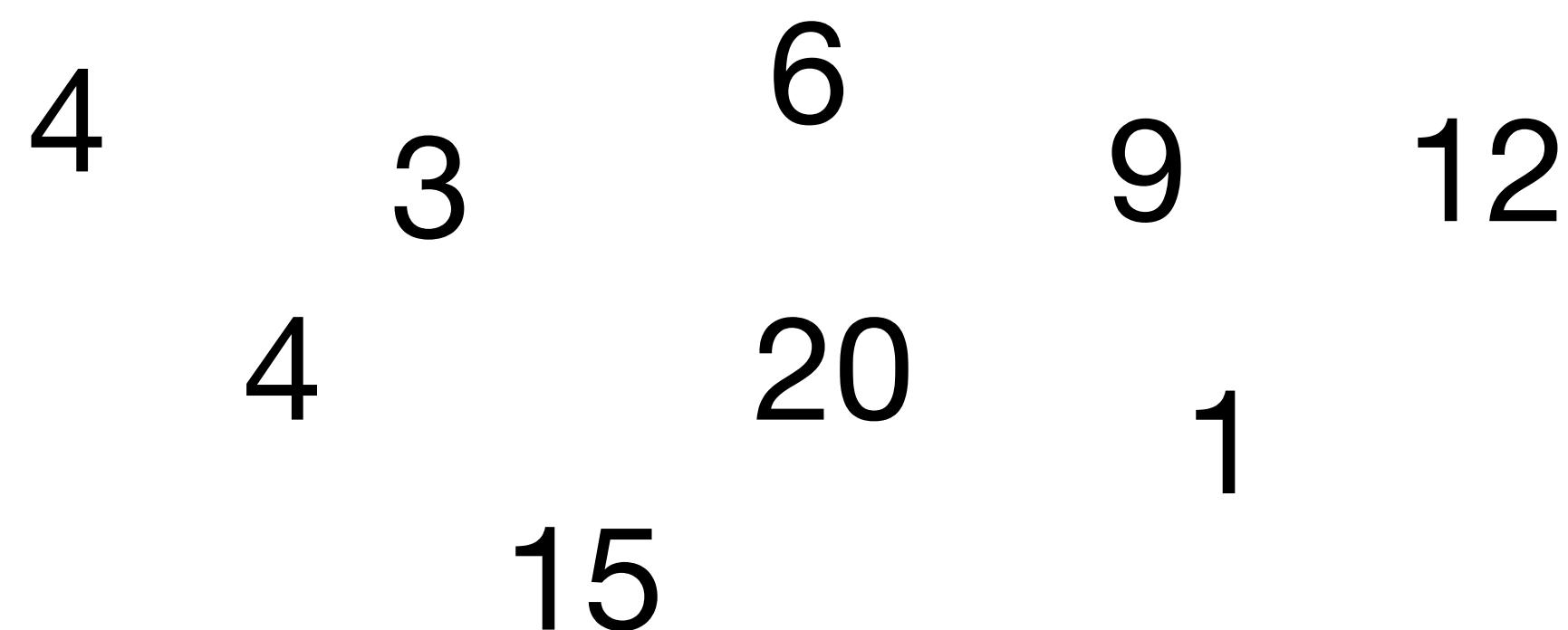
Center



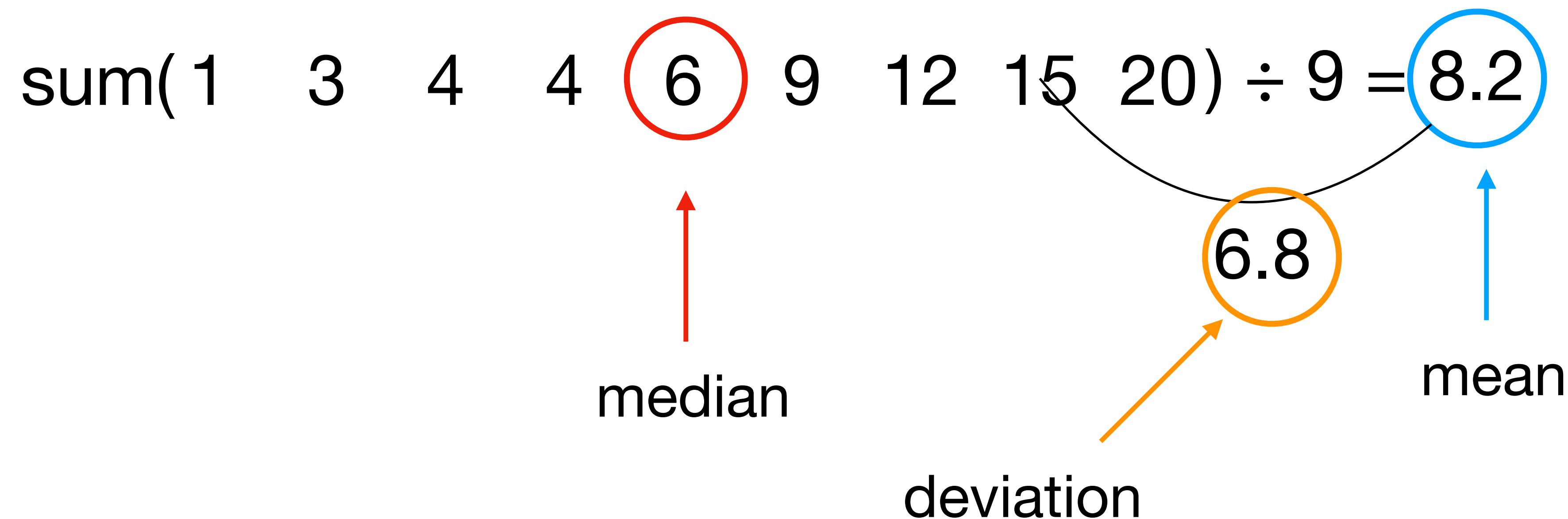
Spread



# Measures of center: median vs. mean



# Measures of center: median vs. mean



- Median is a more **robust** statistic
- $\text{median}(x)$
- Mean can be influenced by extremes
- $\text{mean}(x)$

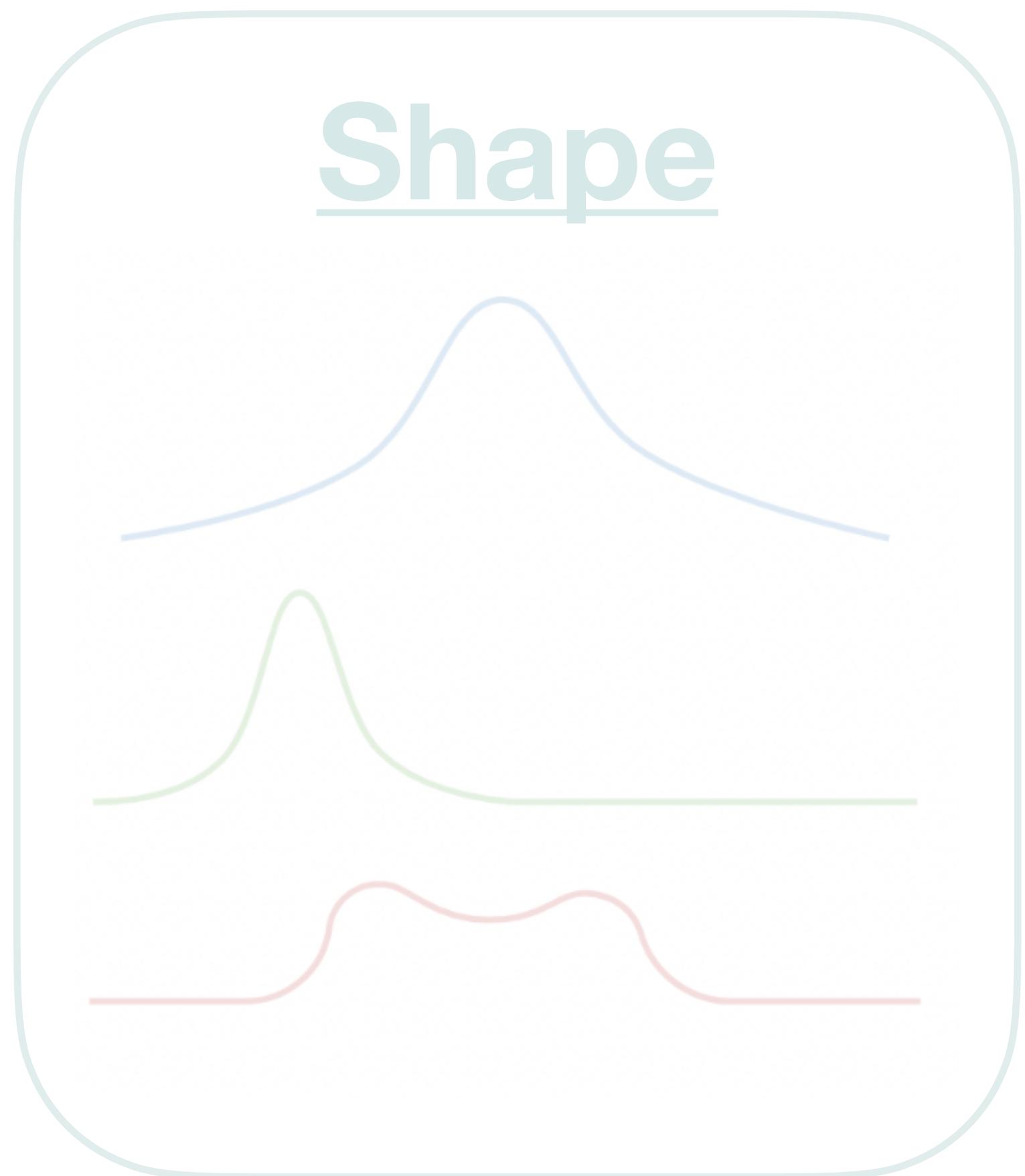
# Measures of center: median vs. mean



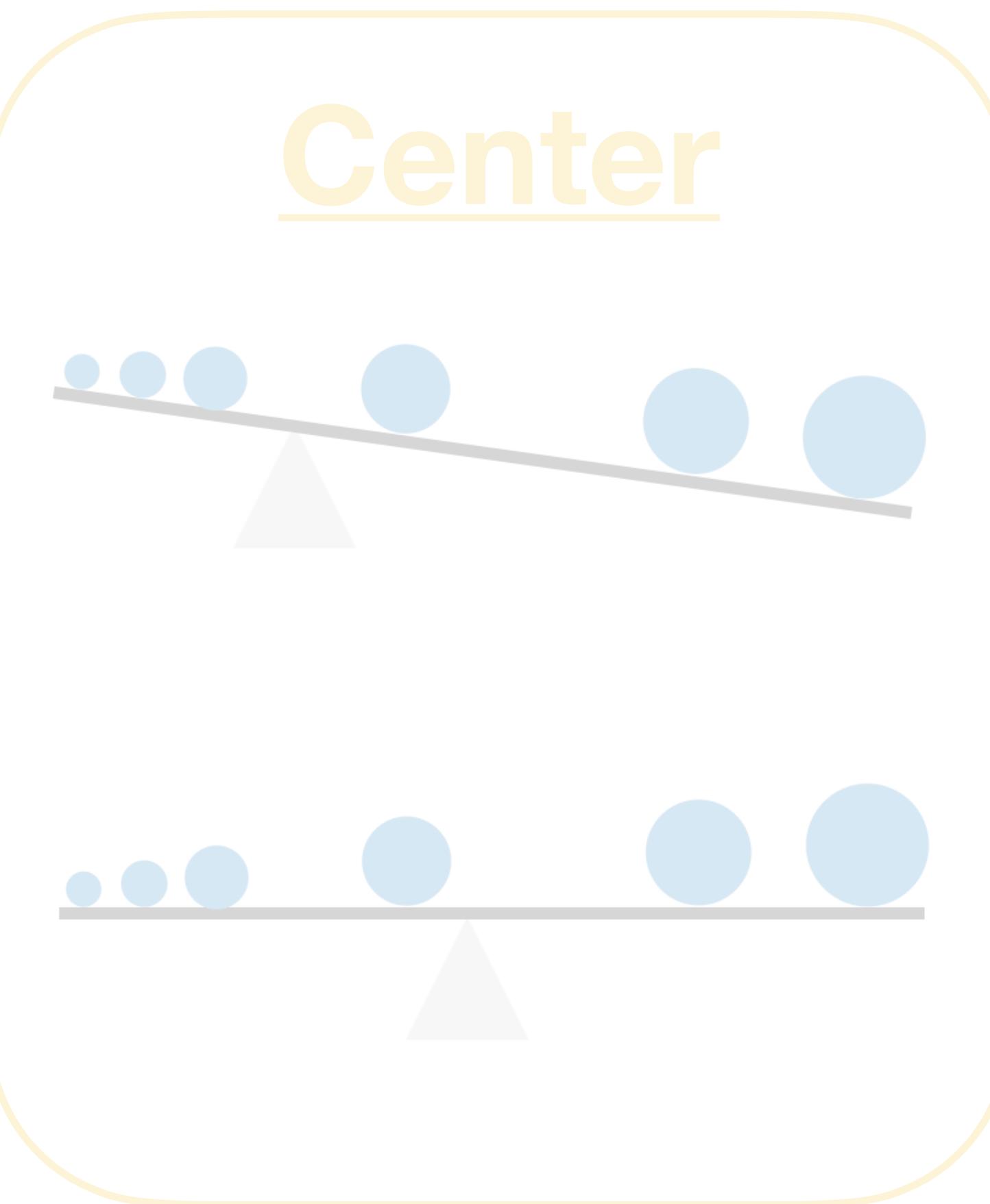
- Median is a more **robust** statistic
- $\text{median}(x)$
- Mean can be influenced by extremes
- $\text{mean}(x)$

# Data exploration and summarization

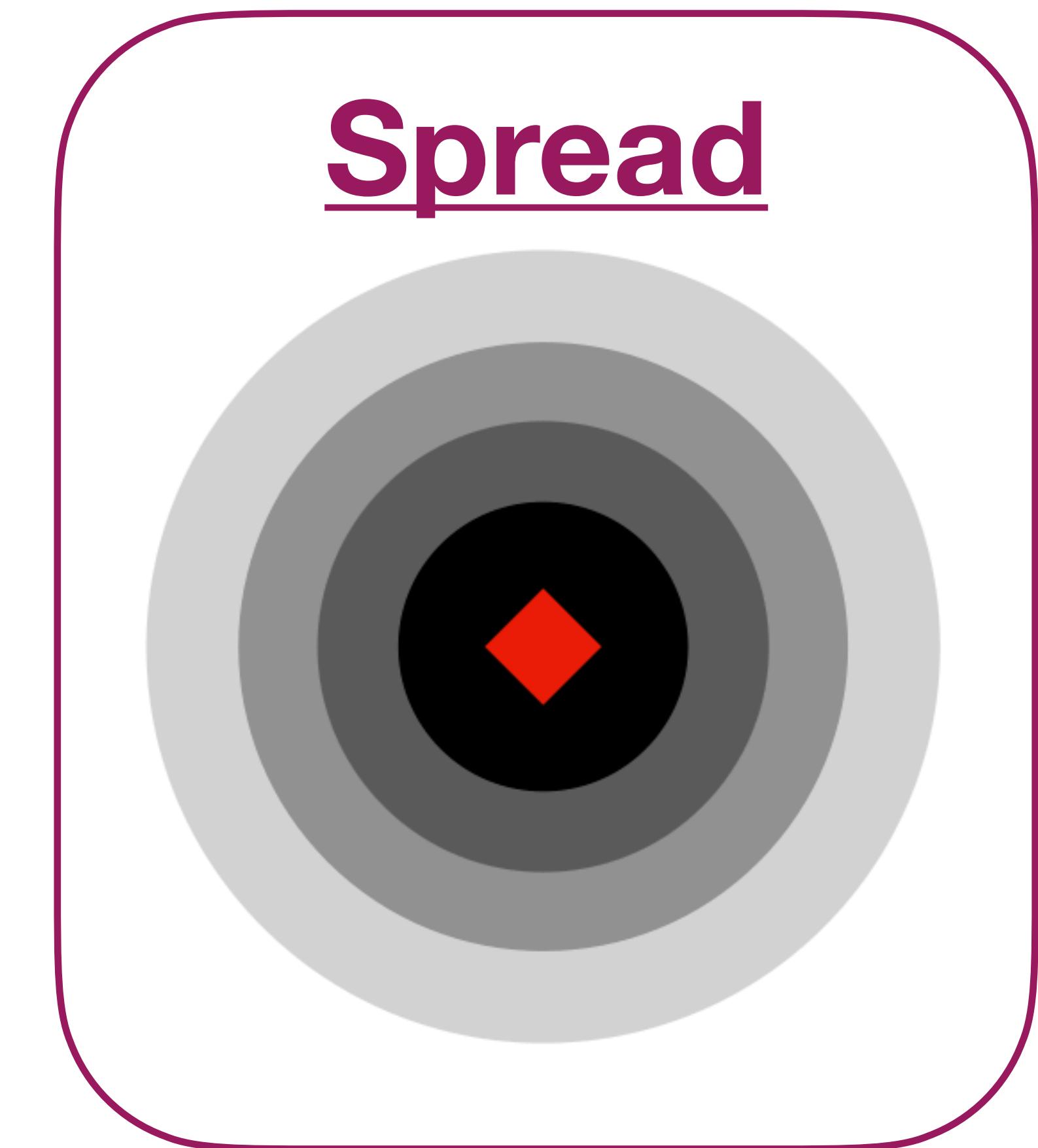
Shape



Center



Spread



# Measures of dispersion: range and IQR

-10, 0, 10, 20, 30

8, 9, 10, 11, 12

**median:** 10

**mean:** 10

???

**median:** 10

**mean:** 10

# Measures of dispersion: range and IQR

-10, 0, 10, 20, 30

40

20

range

IQR

median: 10

mean: 10

8, 9, 10, 11, 12

4

2

median: 10

mean: 10

# Measures of dispersion: range and IQR

-10, 0, 10, 20, 30

40

20

range

IQR

> median: 10

mean: 10

8, 9, 10, 11, 12

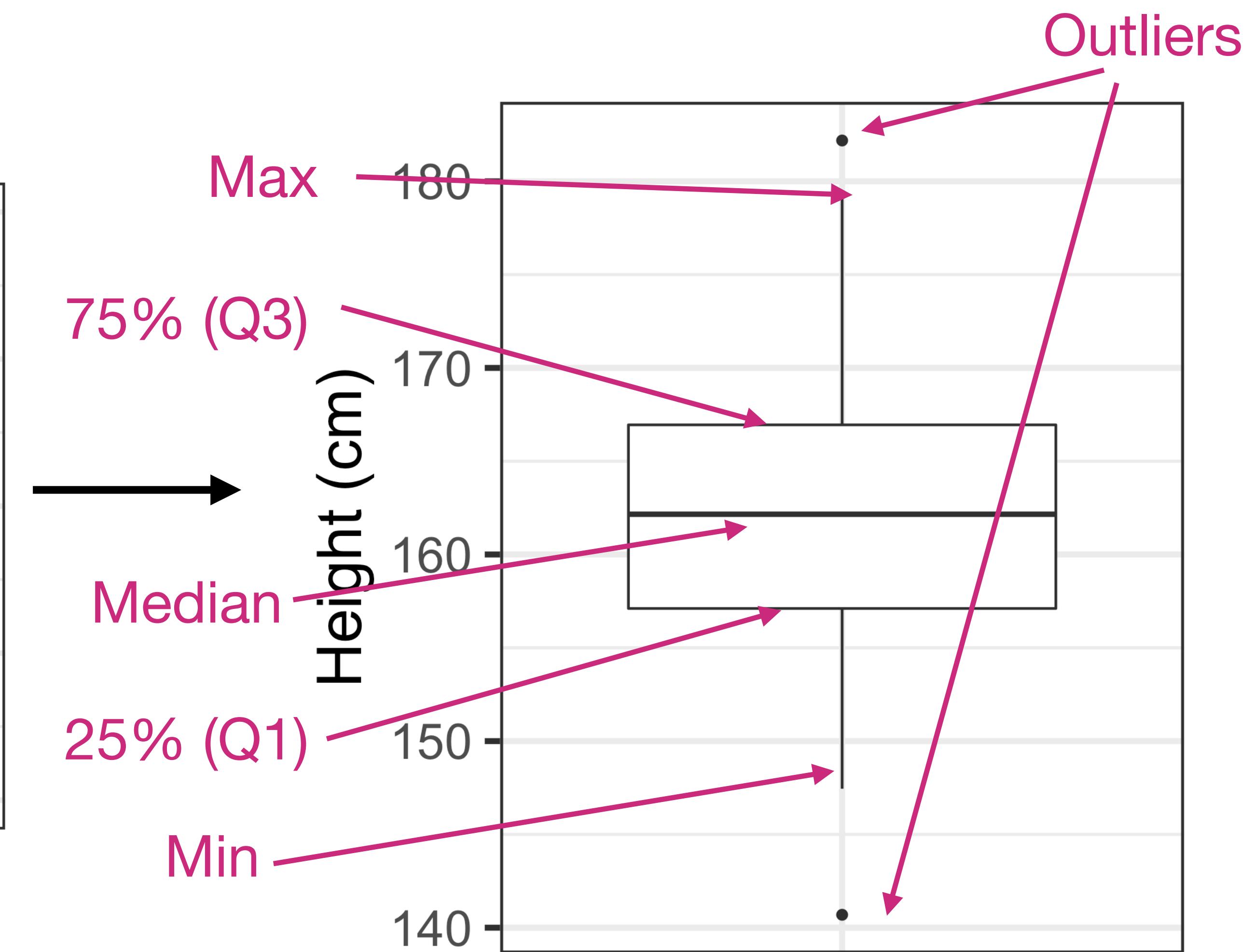
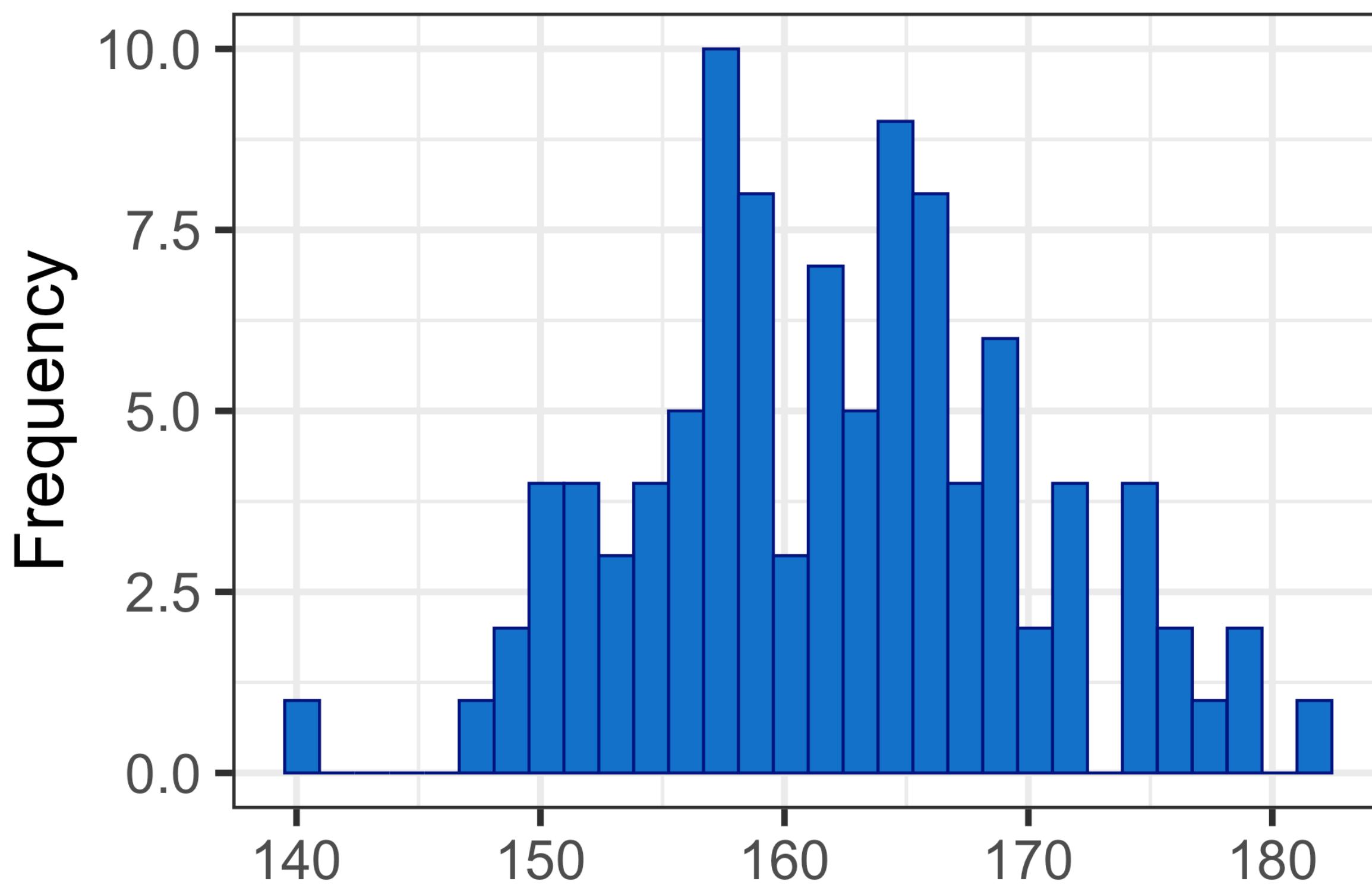
4

2

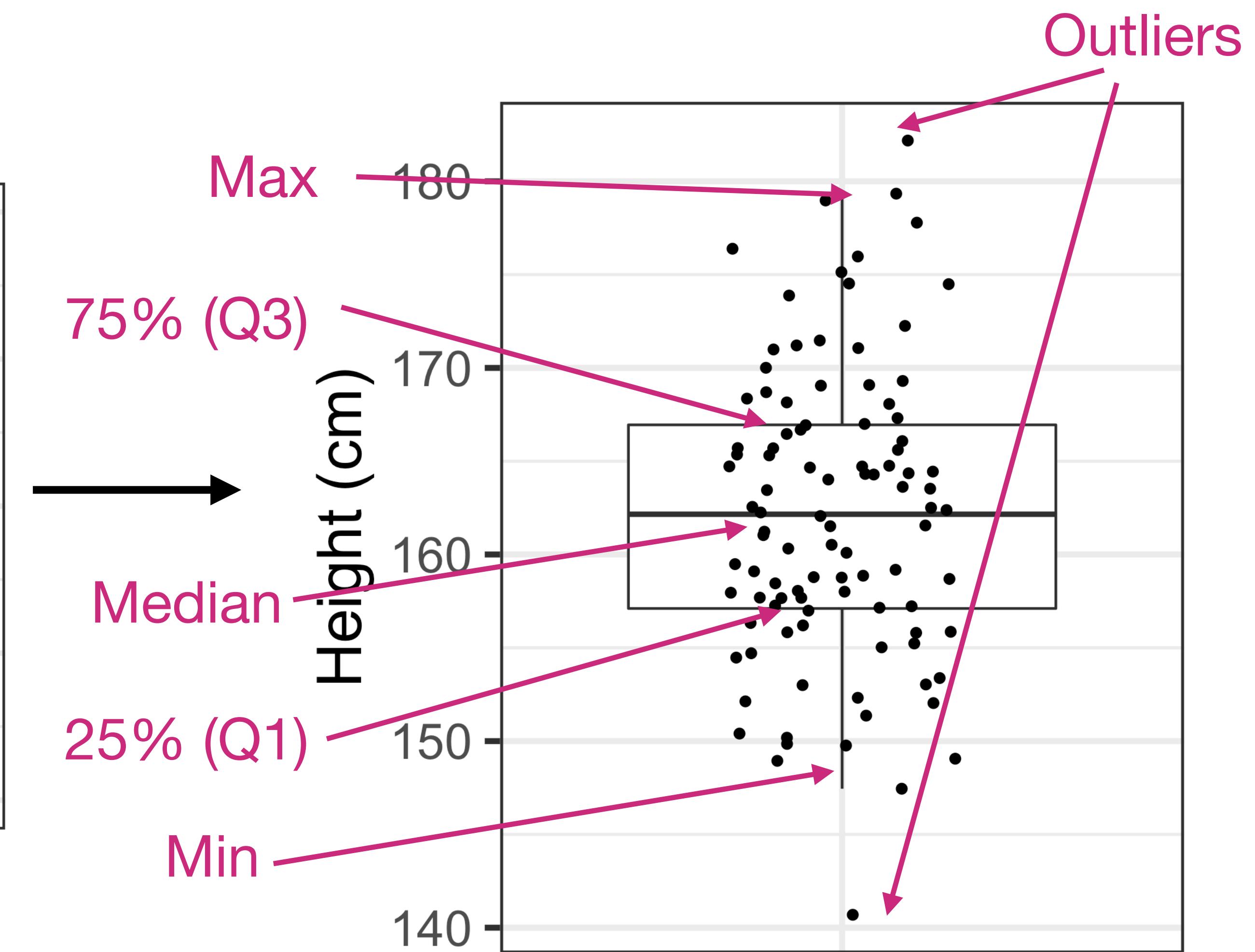
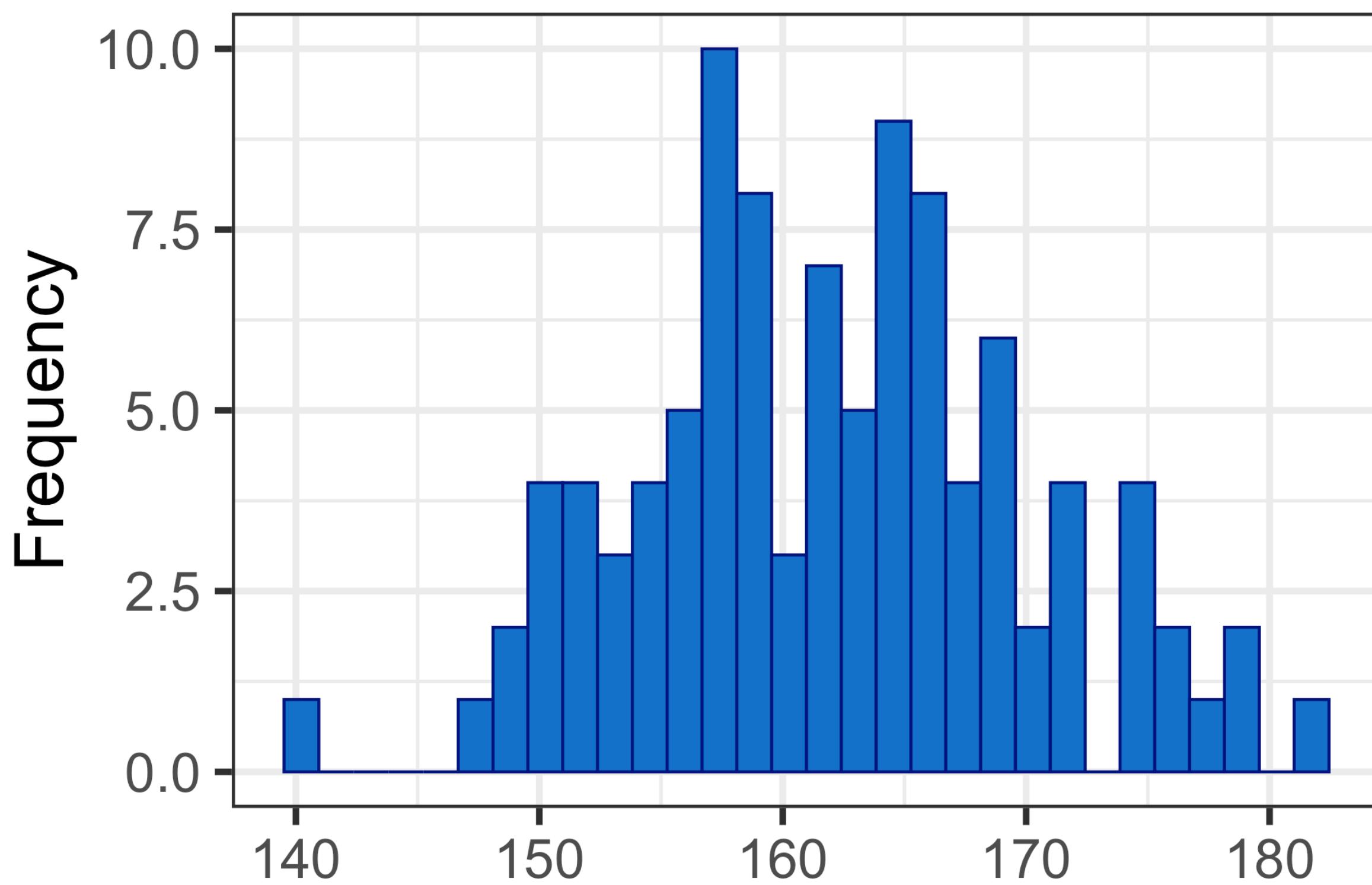
> median: 10

mean: 10

# Plotting distributions with boxplots



# Plotting distributions with boxplots



# Measures of dispersion: standard deviation

-10, 0, 10, 20, 30

8, 9, 10, 11, 12

**Variance:** Mean of the squared deviations

**Standard deviation:** square root of the variance

**median:** 10

> **mean:** 10

**median:** 10

> **mean:** 10

# Measures of dispersion: standard deviation

-10, 0, 10, 20, 30

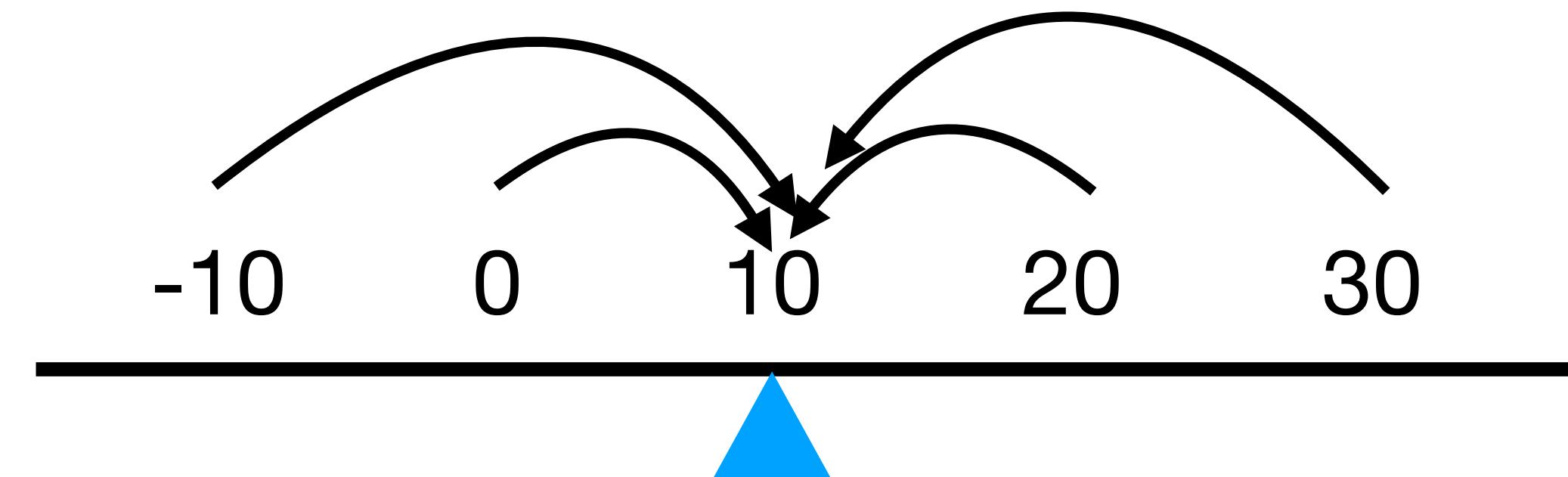
mean: 10

8, 9, 10, 11, 12

20, 10, 0, -10, -20

deviations

2, 1, 0, 1, 2



Variance: Mean of the squared deviations

# Measures of dispersion: standard deviation

-10, 0, 10, 20, 30

mean: 10

8, 9, 10, 11, 12

20, 10, 0, -10, -20

deviations

2, 1, 0, 1, 2

Mean:

400, 100, 0, 100, 400

deviations<sup>2</sup>

Mean:

4, 1, 0, 1, 4

variance

200

2

Variance: Mean of the squared deviations

# Measures of dispersion: standard deviation

-10, 0, 10, 20, 30

mean: 10

8, 9, 10, 11, 12

20, 10, 0, -10, -20

deviations

2, 1, 0, 1, 2

400, 100, 0, 100, 400

deviations<sup>2</sup>

4, 1, 0, 1, 4

200

variance

2

14.1

standard deviation

1.41

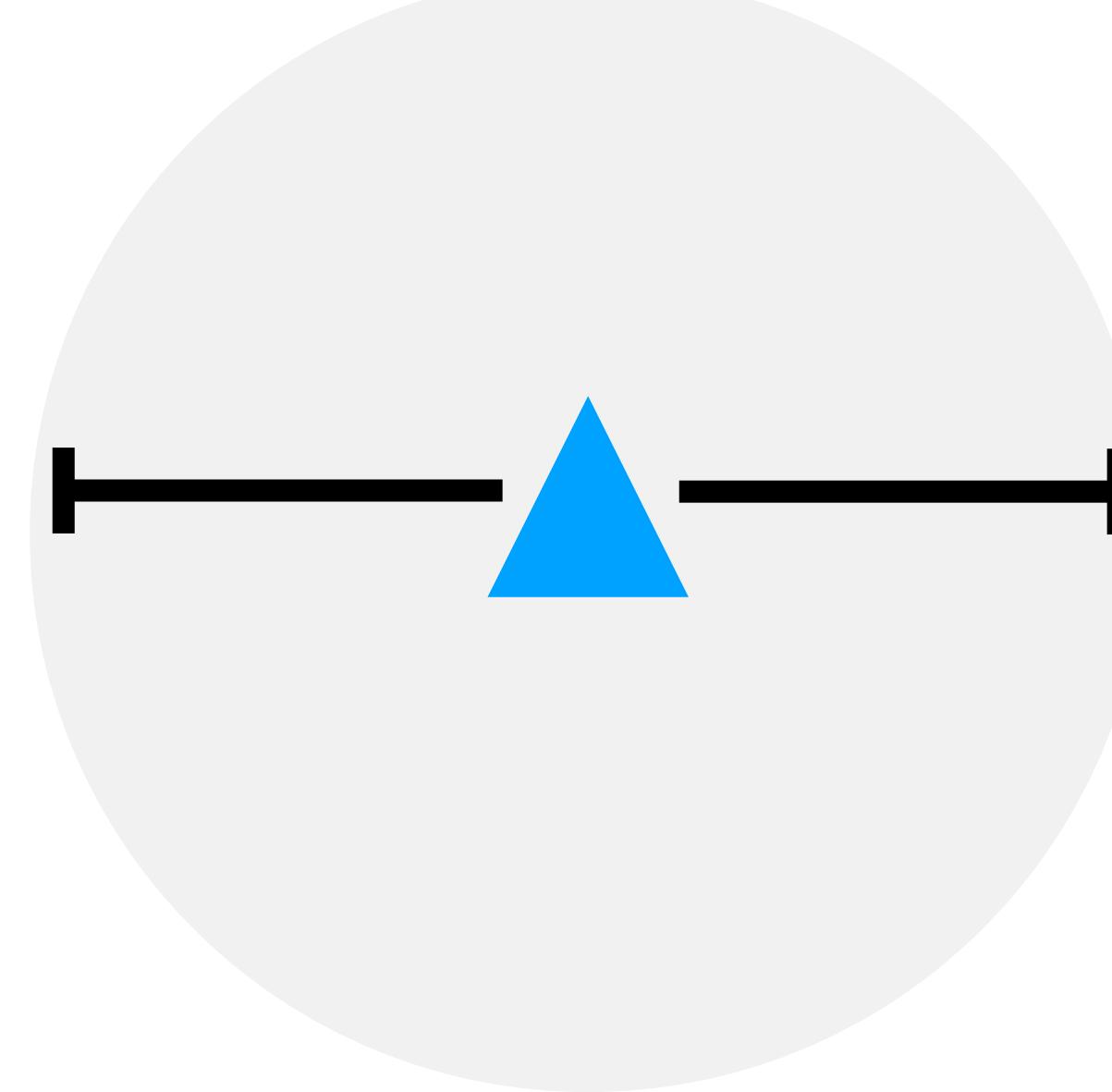
**Standard deviation:** square root of the variance

# Measures of dispersion: standard deviation

-10, 0, 10, 20, 30

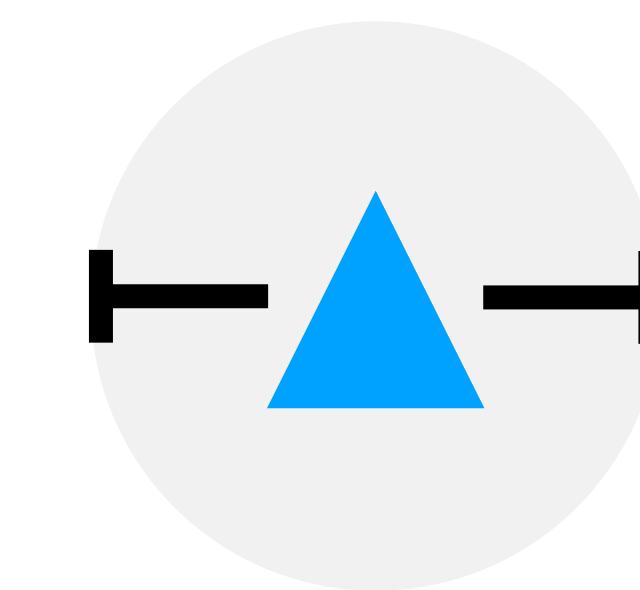
mean: 10

8, 9, 10, 11, 12



14.1

standard deviation



1.41

*“typical distance  
of the observations  
from their mean”*



**Standard deviation:** square root of the variance

# Measures of dispersion: standard deviation

> `var(x)`

> `sd(x)`

$\sigma =$

$$\sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

The diagram shows the formula for standard deviation. A teal rounded rectangle surrounds the entire formula. Inside, a pink rounded rectangle highlights the numerator  $\sum (x_i - \mu)^2$ . Within this pink area, an orange rounded rectangle highlights the term  $(x_i - \mu)^2$ . Above the formula, the word "deviations" is written in orange. To the right of the formula, four labels are aligned vertically: "Squared deviations" (grey), "variance" (pink), and "Standard deviation" (teal). The label "Standard deviation" is positioned below the teal rounded rectangle.

$\sigma$ : population standard deviation

$N$ : Size of population

$x_i$ : observed values

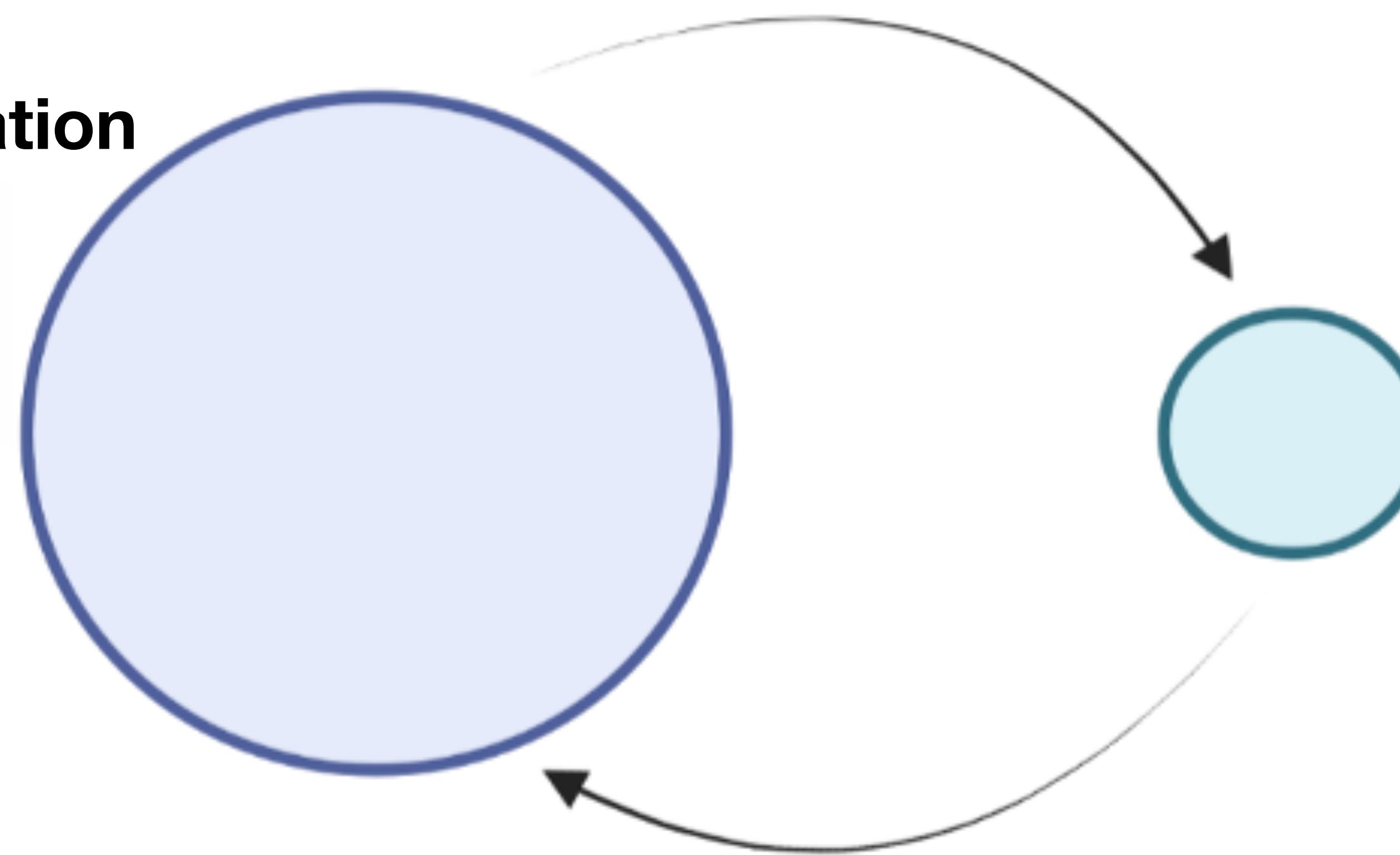
$\mu$ : population mean

Similar, but different, equations describe a population vs. a sample

*Random sampling*

**Population**

$$\sigma = \sqrt{\frac{\sum(x_i - \mu)^2}{N}}$$



**Sample**

$$s = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1}}$$

*Inference*

# Measures of dispersion: standard deviation

```
> var(x)
```

```
> sd(x)
```

$$s = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1}}$$

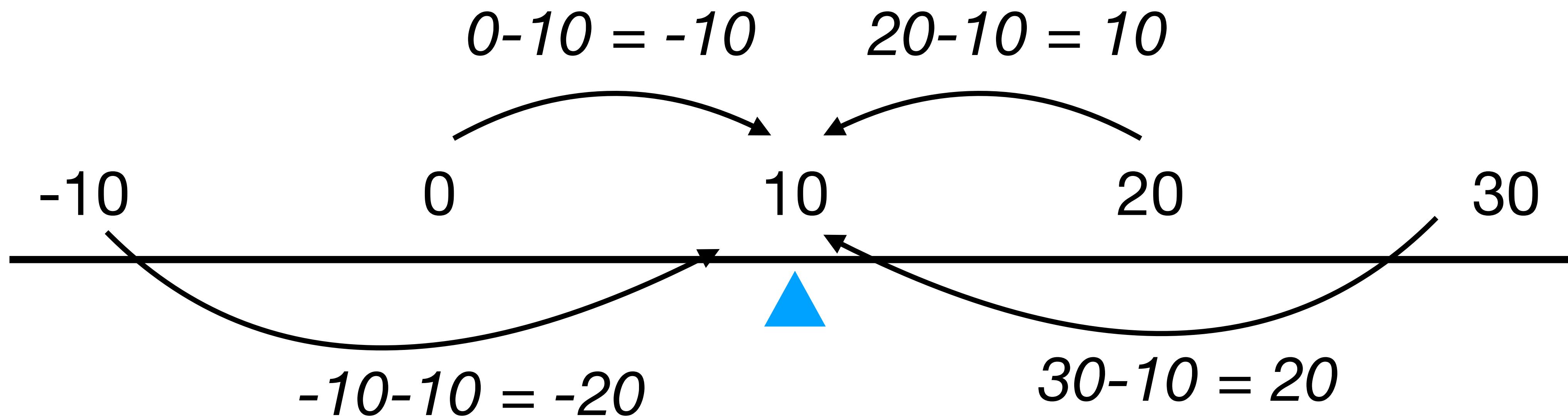
**s:** sample standard deviation

**N:** number of observations

$x_i$ : observed values

$\bar{x}$ : mean value of samples

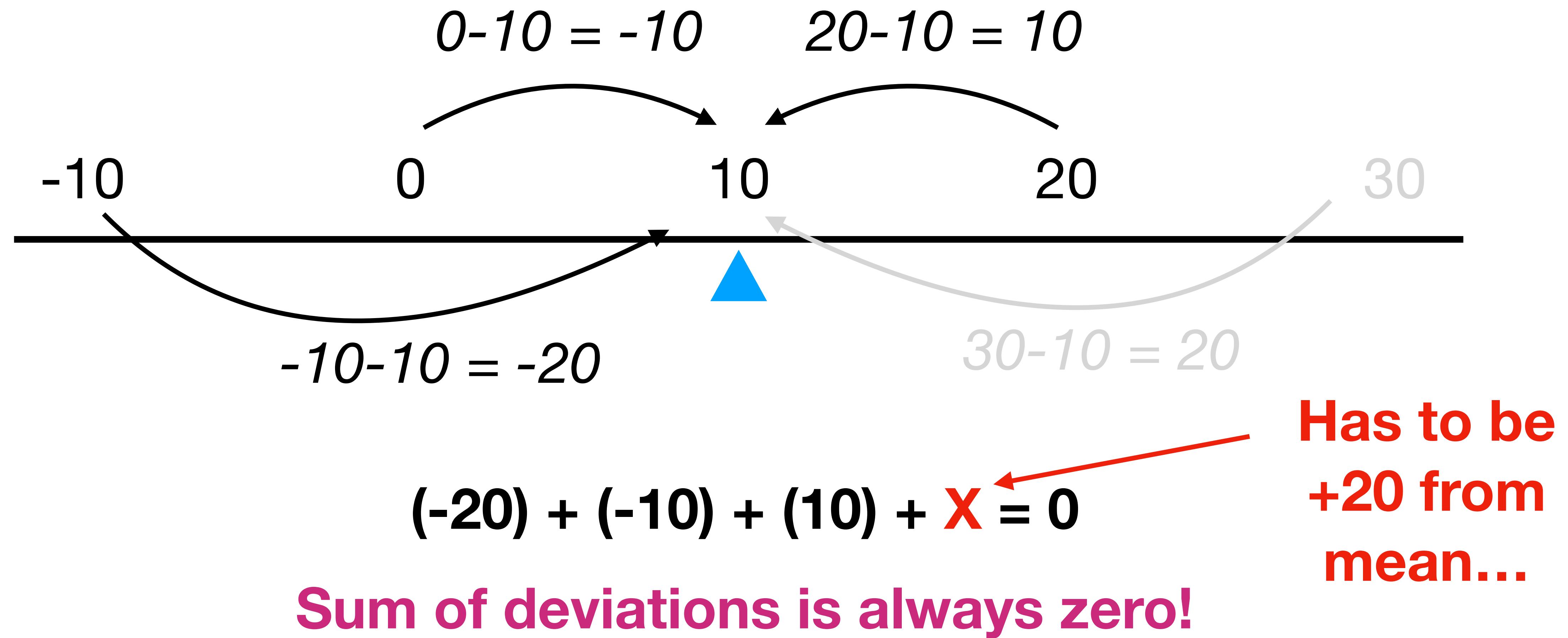
# N-1: degrees of freedom explained



$$(-20) + (-10) + (10) + (20) = 0$$

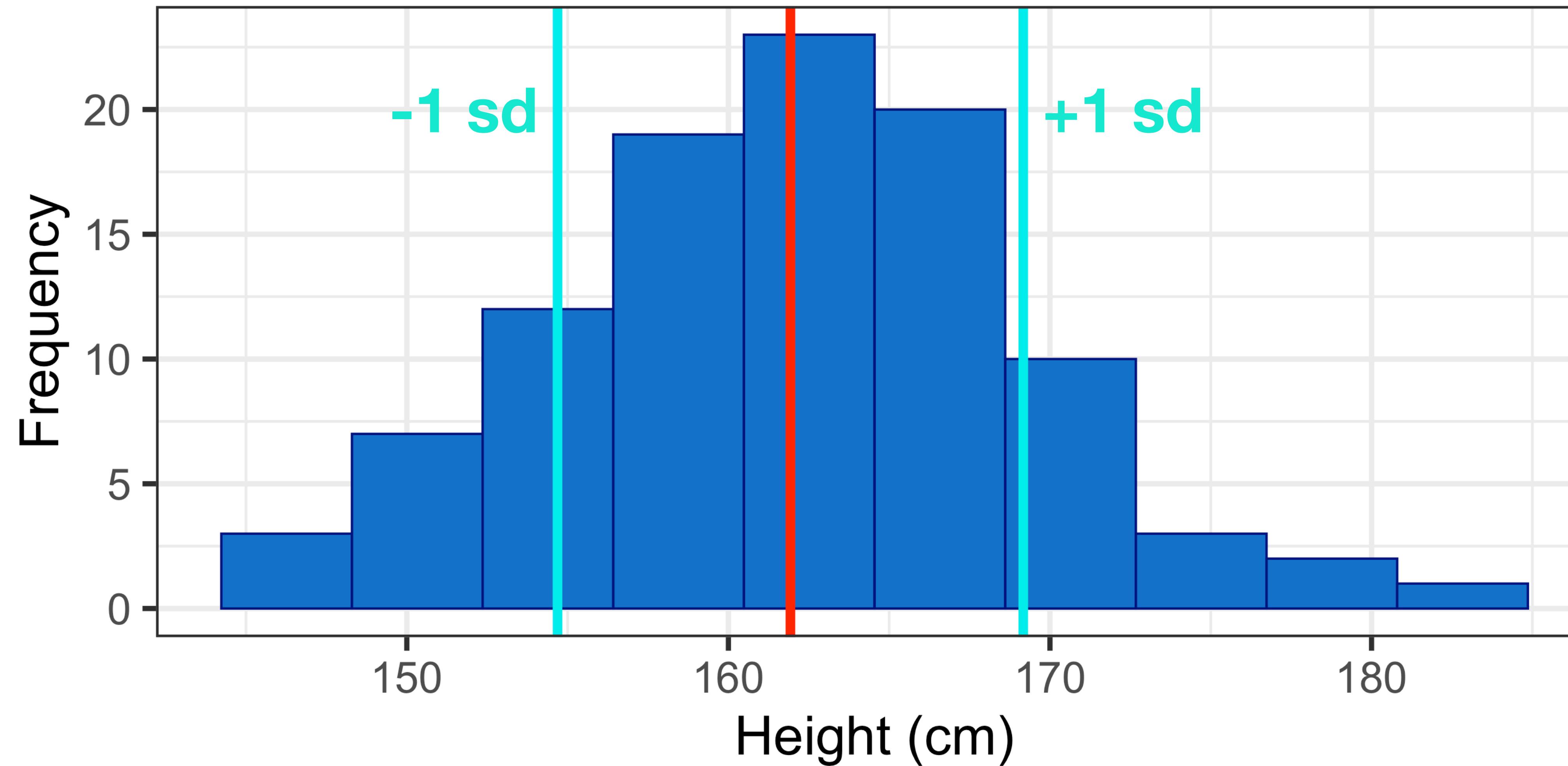
**Sum of deviations is always zero!**

# N-1: degrees of freedom explained

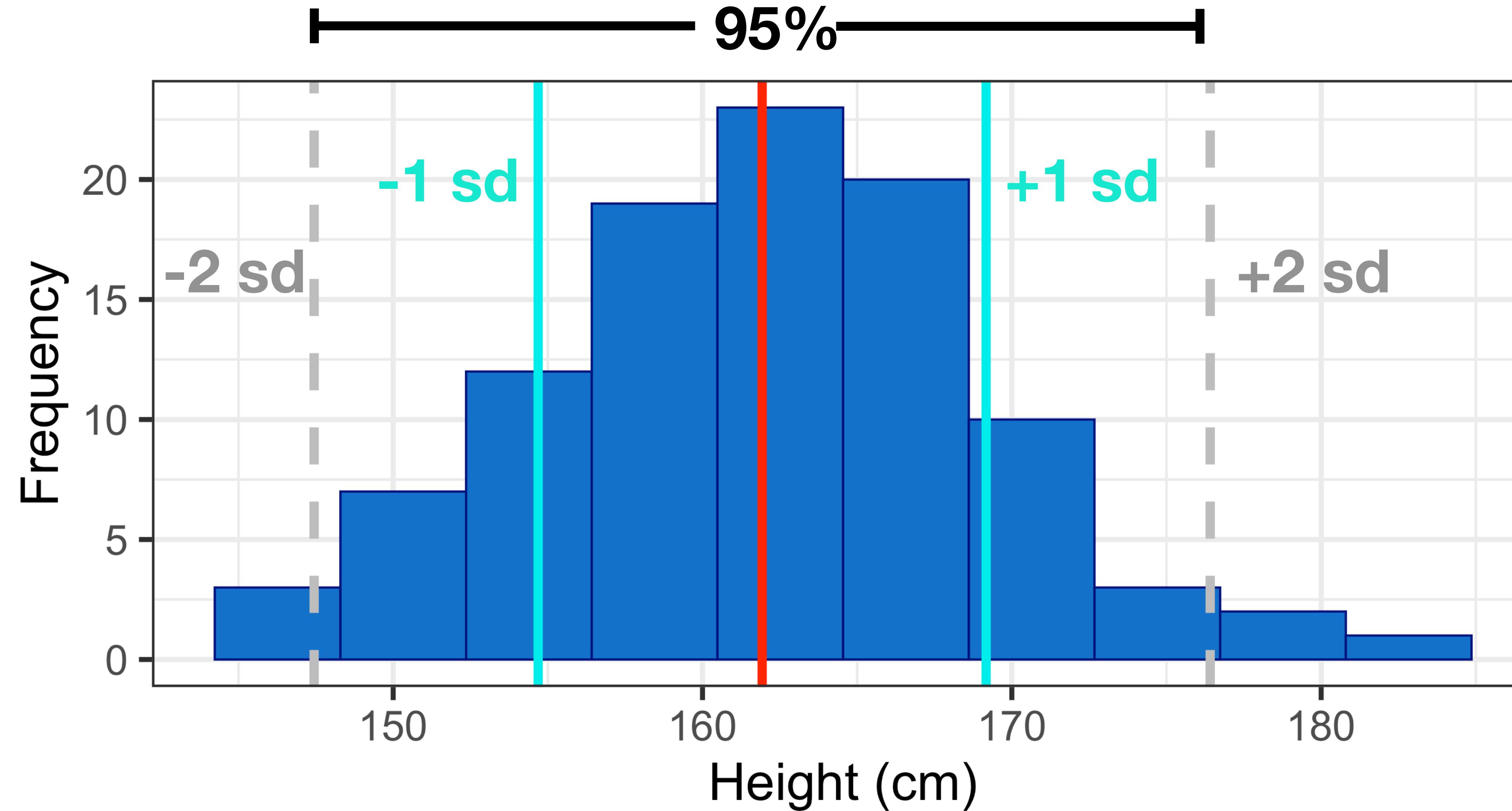


# Measures of dispersion: standard deviation

68%



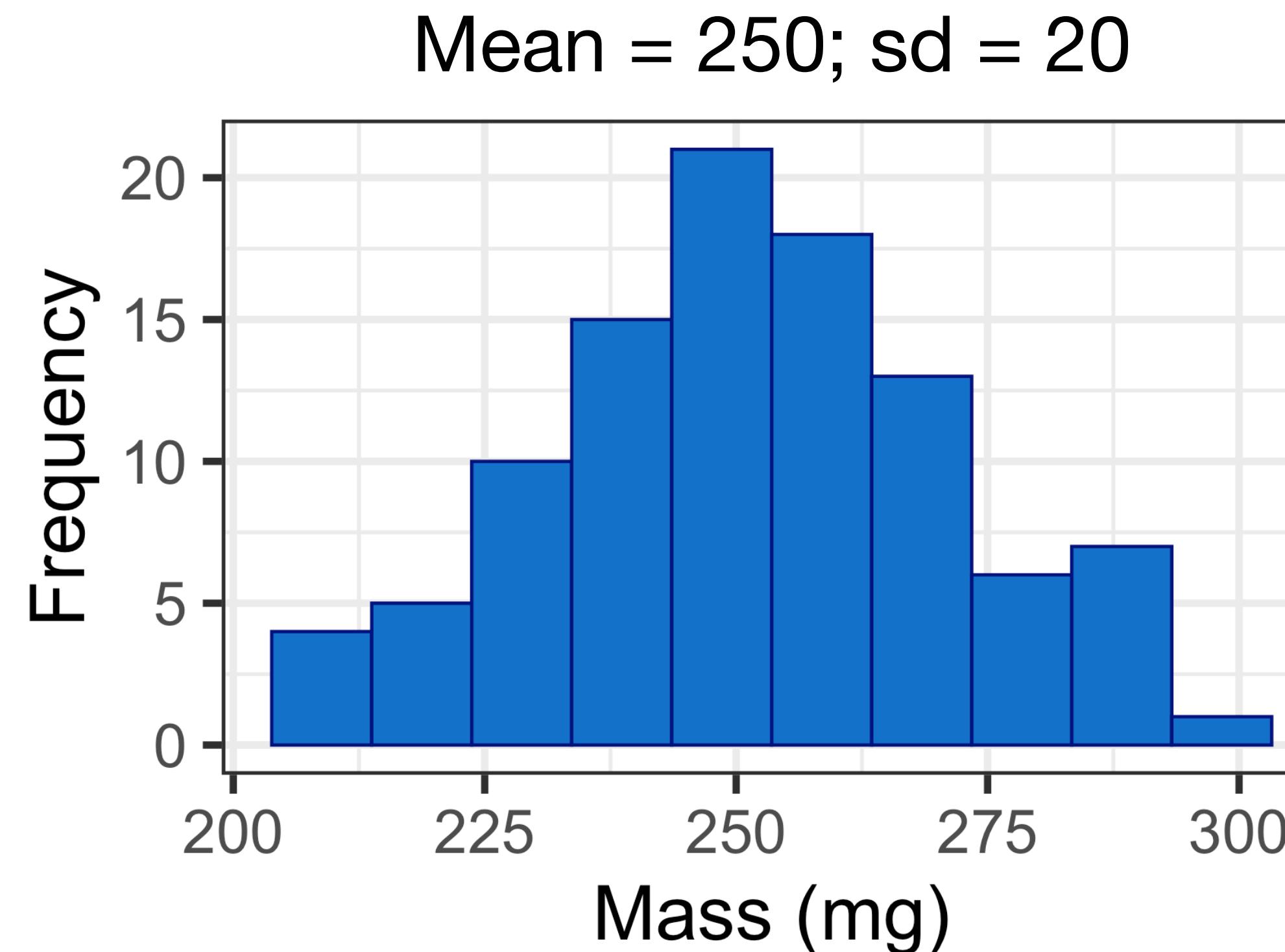
# Measures of dispersion: standard deviation



# Linear transformation of variables

**Suppose you collected data in milligrams but need to convert to grams.**

```
new_data = (1/1000) * old_data
```

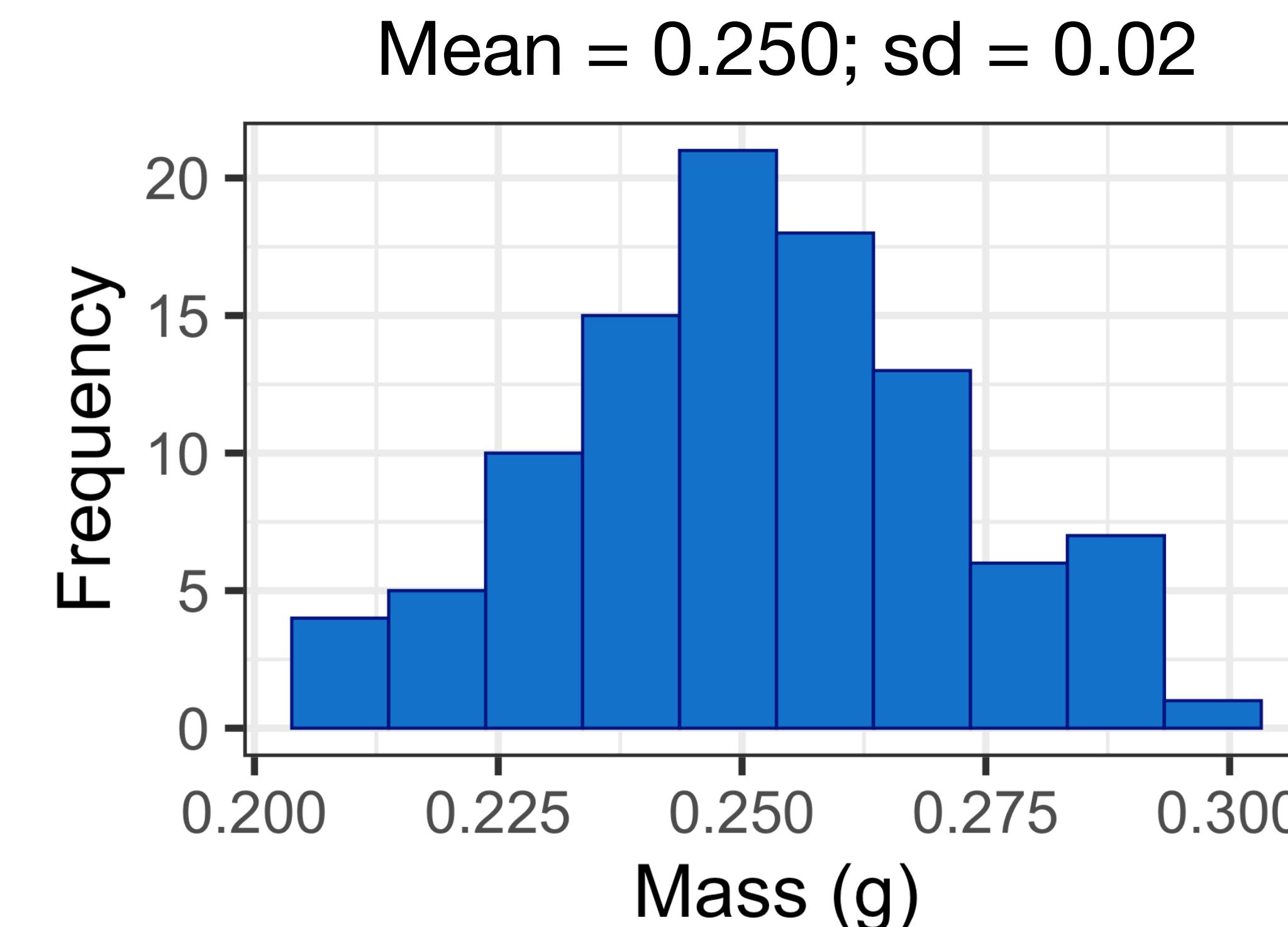
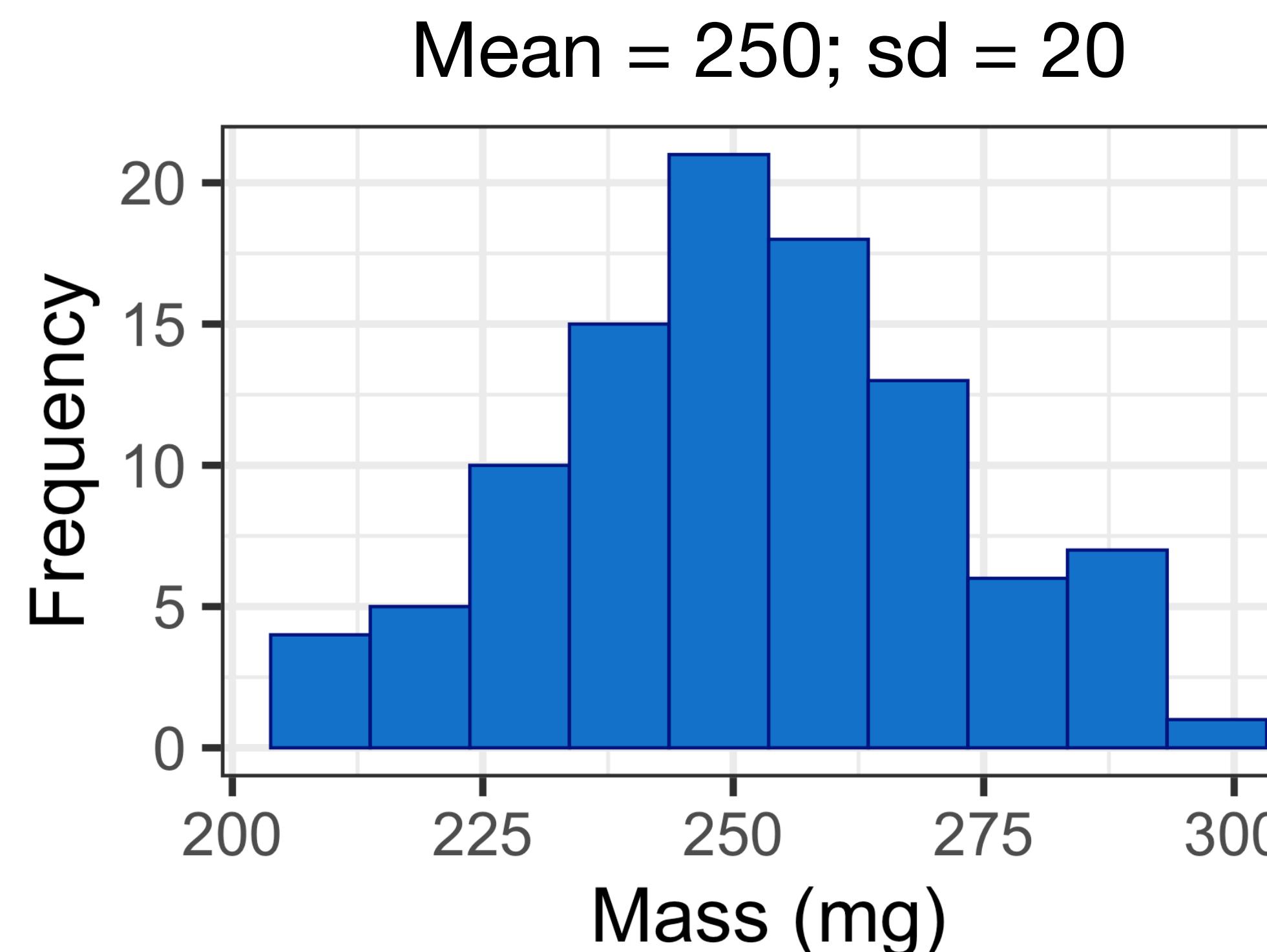


# Linear transformation of variables

**Suppose you collected data in milligrams but need to convert to grams.**

```
new_data = (1/1000) * old_data
```

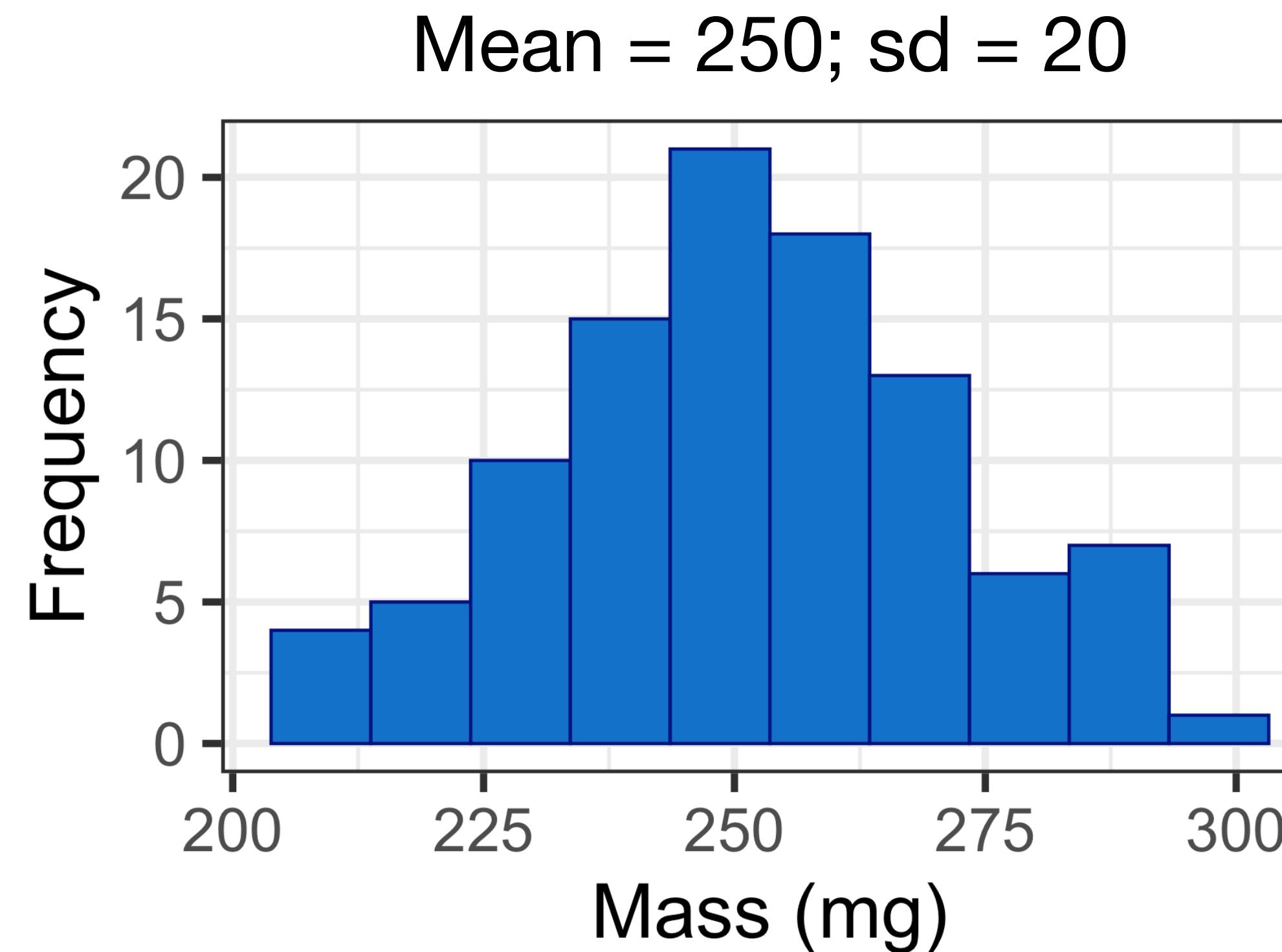
**Mean & SD scale**



# Linear transformation of variables

**Suppose your scale is off by 50 mg...**

$$\text{new\_data} = \text{old\_data} + 50$$

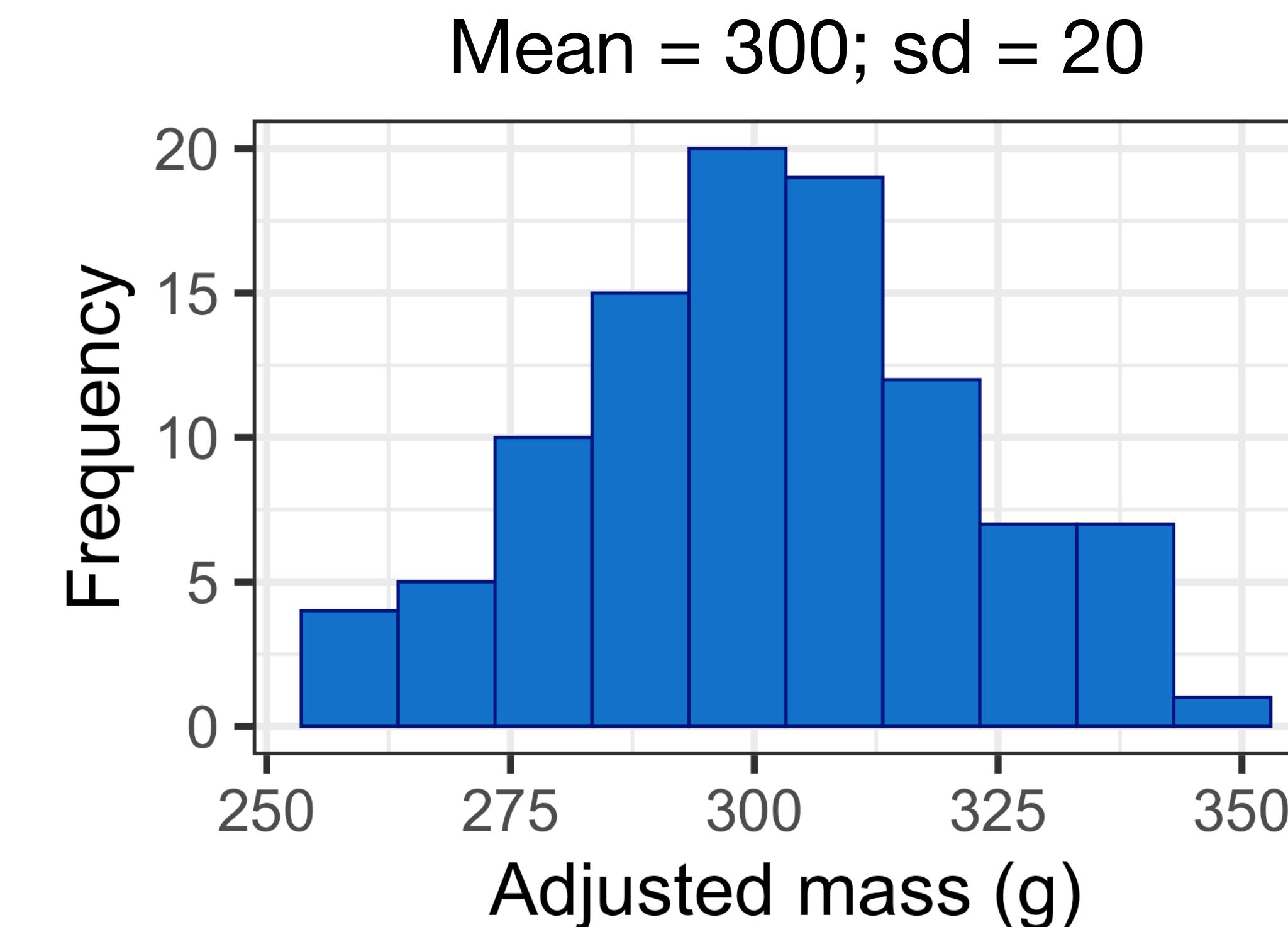
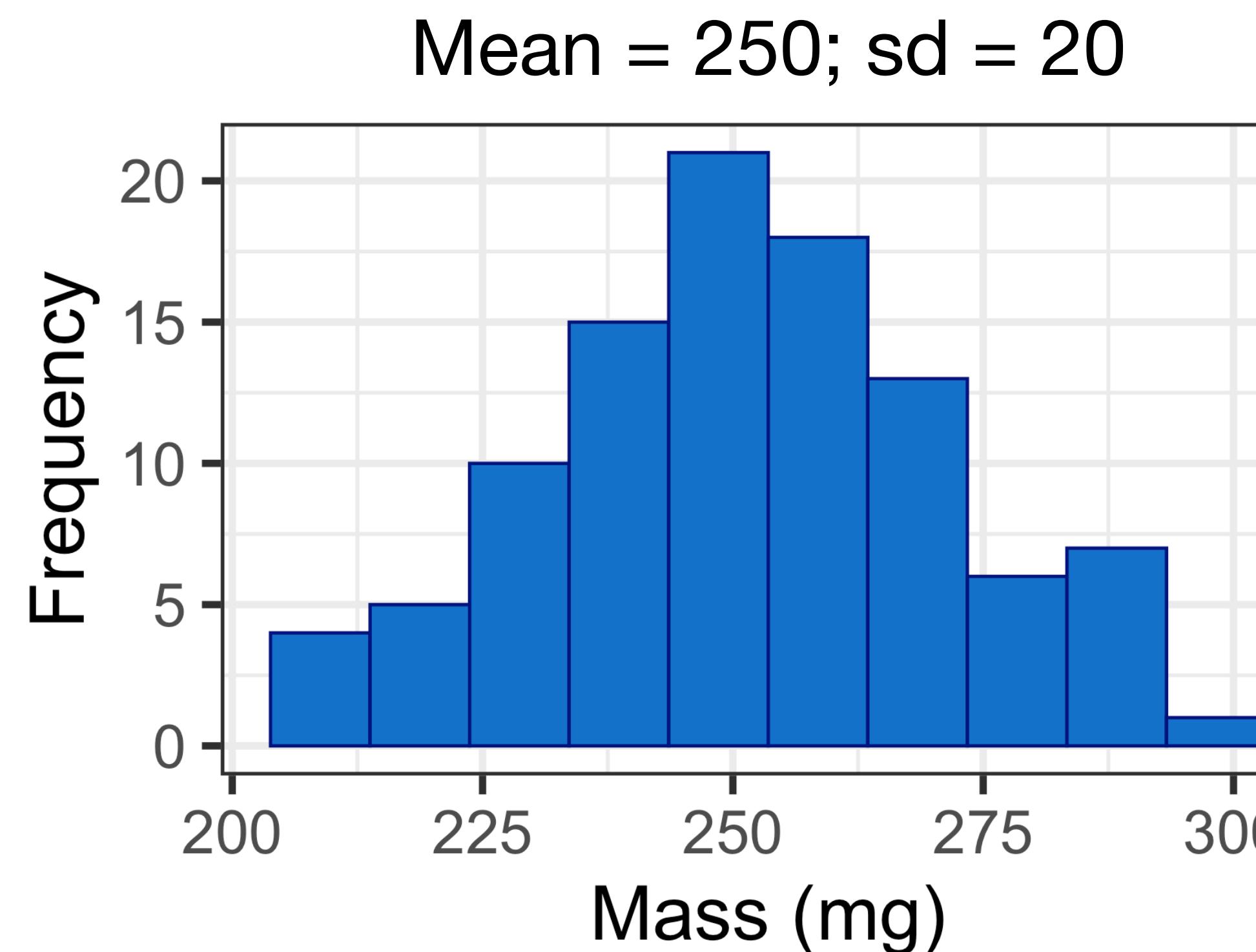


# Linear transformation of variables

**Suppose your scale is off by 50 mg...**

`new_data = old_data + 50`

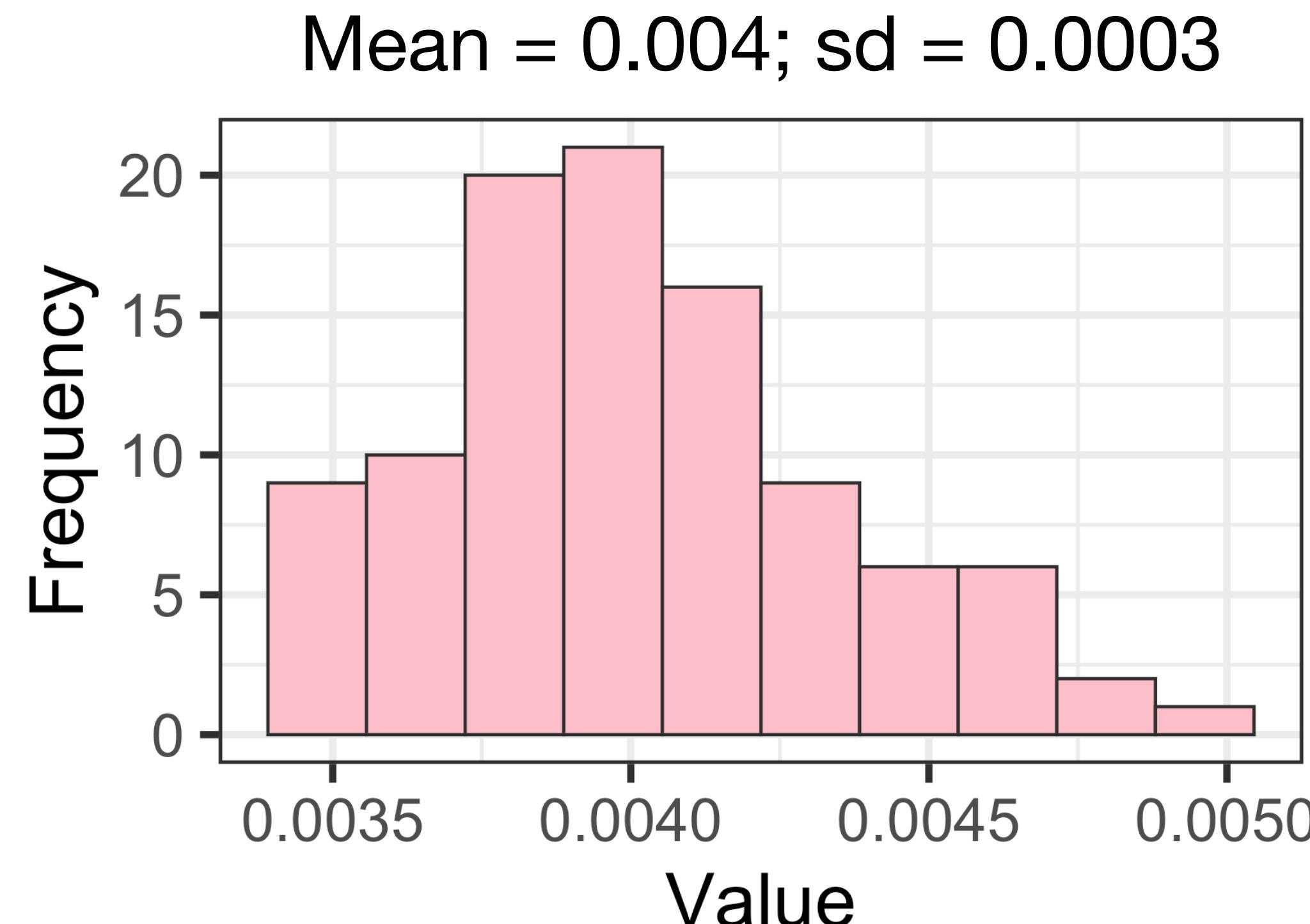
**Mean scales, SD not**



# Non-linear transformation of variables

**Suppose your data has a strong right skew...**

```
new_data = old_data + 50
```



# Non-linear transformation of variables

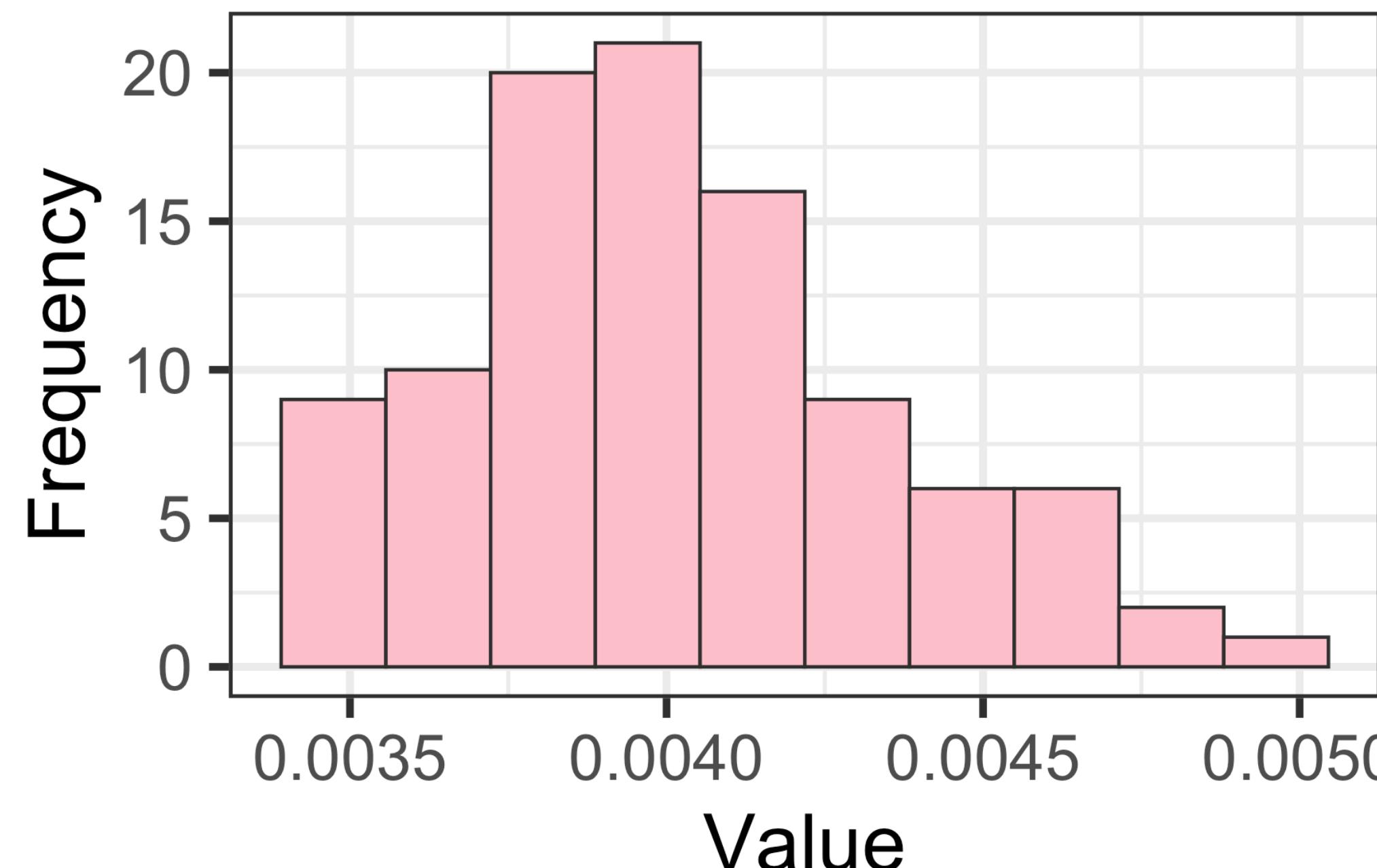
- $Y' = \sqrt{Y}$
- $Y' = \log(Y)$
- $Y' = 1/Y$
- $Y' = Y^2$

Suppose your data has a strong right skew...

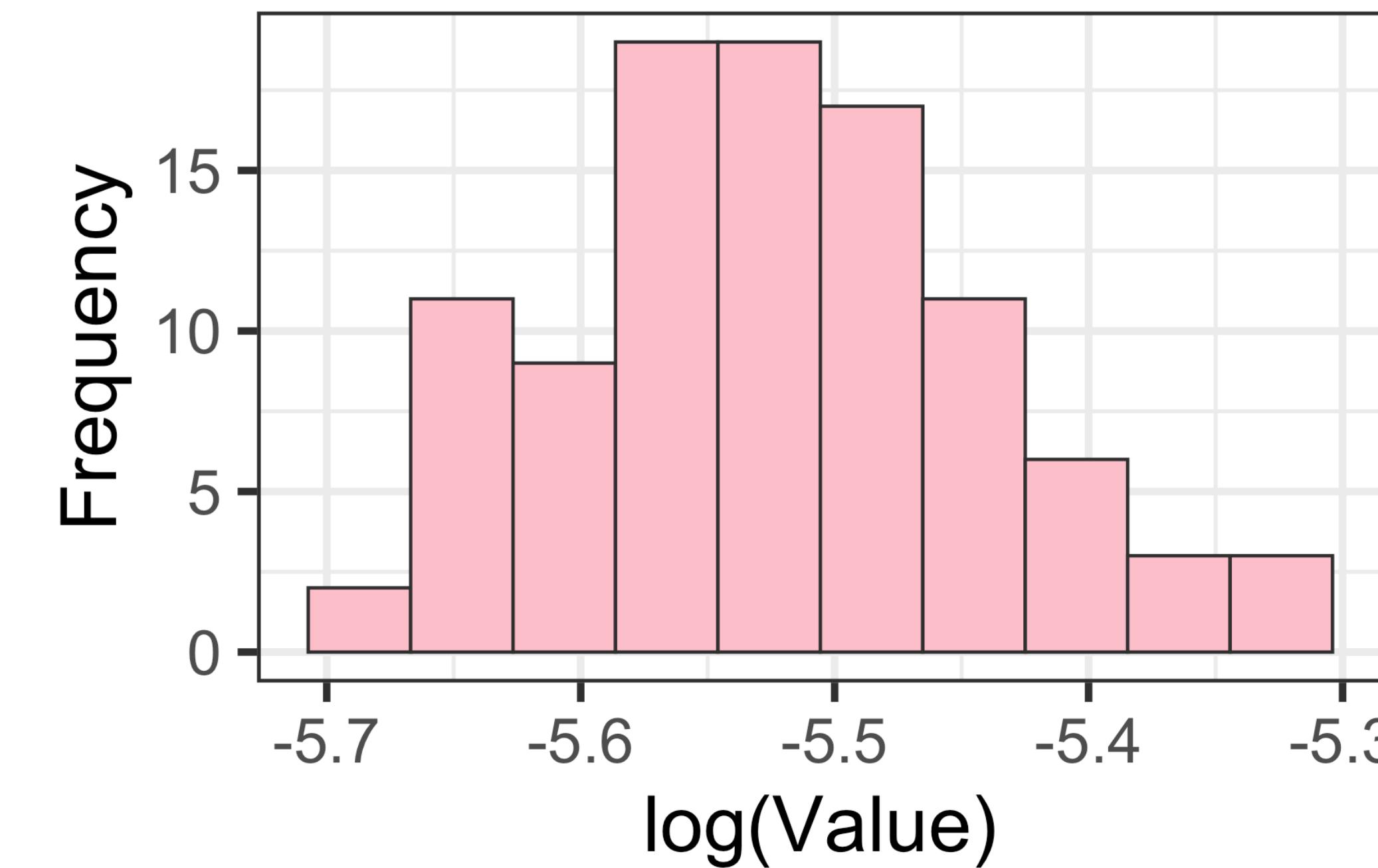
`new_data = old_data + 50`

**Mean and SD scale non-linearly**

Mean = 0.004; sd = 0.0003



Mean = -5.5; sd = 0.08



# Data Visualization with ggplot2 :: CHEAT SHEET

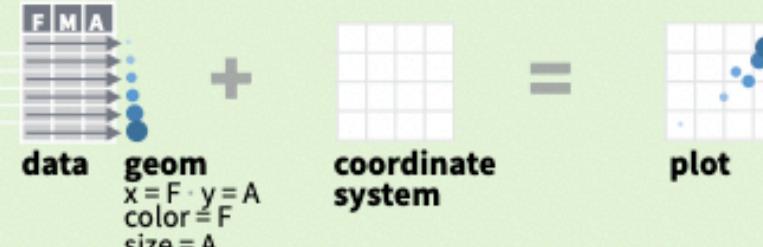


## Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and geoms—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot (data = <DATA>) +
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),
  stat = <STAT>, position = <POSITION>) +
  <COORDINATE_FUNCTION>+
  <FACET_FUNCTION>+
  <SCALE_FUNCTION>+
  <THEME_FUNCTION>
```

[ required ]

Not required, sensible defaults supplied

`ggplot(data = mpg, aes(x = cty, y = hwy))` Begins a plot that you finish by adding layers to. Add one geom function per layer.

**aesthetic mappings** **data** **geom**

`qplot(x = cty, y = hwy, data = mpg, geom = "point")` Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

`last_plot()` Returns the last plot

`ggsave("plot.png", width = 5, height = 5)` Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.



## Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
```

- a + geom\_blank()**  
(Useful for expanding limits)
- b + geom\_curve(aes(yend = lat + 1, xend = long + 1, curvature = z))** - x, xend, y, yend, alpha, angle, color, curvature, linetype, size
- a + geom\_path(lineend = "butt", linejoin = "round", linemitre = 1)**  
x, y, alpha, color, group, linetype, size
- a + geom\_polygon(aes(group = group))**  
x, y, alpha, color, fill, group, linetype, size
- b + geom\_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1))** - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size
- a + geom\_ribbon(aes(ymax = unemploy - 900, ymin = unemploy + 900))** - x, ymax, ymin, alpha, color, fill, group, linetype, size

### LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

- b + geom\_abline(aes(intercept = 0, slope = 1))**
- b + geom\_hline(aes(yintercept = lat))**
- b + geom\_vline(aes(xintercept = long))**
- b + geom\_segment(aes(yend = lat + 1, xend = long + 1))**
- b + geom\_spoke(aes(angle = 1:1155, radius = 1))**

### ONE VARIABLE continuous

- ```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```
- c + geom\_area(stat = "bin")**  
x, y, alpha, color, fill, linetype, size
  - c + geom\_density(kernel = "gaussian")**  
x, y, alpha, color, fill, group, linetype, size, weight
  - c + geom\_dotplot()**  
x, y, alpha, color, fill
  - c + geom\_freqpoly()** x, y, alpha, color, group, linetype, size
  - c + geom\_histogram(binwidth = 5)** x, y, alpha, color, fill, linetype, size, weight
  - c2 + geom\_qq(aes(sample = hwy))** x, y, alpha, color, fill, linetype, size, weight

### discrete

- ```
d <- ggplot(mpg, aes(f1))
```
- d + geom\_bar()**  
x, alpha, color, fill, linetype, size, weight

### TWO VARIABLES

#### continuous x , continuous y

- ```
e <- ggplot(mpg, aes(cty, hwy))
```
- e + geom\_label(aes(label = cty, nudge\_x = 1, nudge\_y = 1, check\_overlap = TRUE))** x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

- e + geom\_jitter(height = 2, width = 2)**  
x, y, alpha, color, fill, shape, size

- e + geom\_point()**, x, y, alpha, color, fill, shape, size, stroke

- e + geom\_quantile()**, x, y, alpha, color, group, linetype, size, weight

- e + geom\_rug(sides = "bl")**, x, y, alpha, color, linetype, size

- e + geom\_smooth(method = lm)**, x, y, alpha, color, fill, group, linetype, size, weight

- e + geom\_text(aes(label = cty), nudge\_x = 1, nudge\_y = 1, check\_overlap = TRUE)**, x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

#### discrete x , continuous y

- ```
f <- ggplot(mpg, aes(class, hwy))
```

- f + geom\_col()**, x, y, alpha, color, fill, group, linetype, size
- f + geom\_boxplot()**, x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight
- f + geom\_dotplot(binaxis = "y", stackdir = "center")**, x, y, alpha, color, fill, group
- f + geom\_violin(scale = "area")**, x, y, alpha, color, fill, group, linetype, size, weight

#### discrete x , discrete y

- ```
g <- ggplot(diamonds, aes(cut, color))
```

- g + geom\_count()**, x, y, alpha, color, fill, shape, size, stroke

### THREE VARIABLES

- ```
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2))
```

- ```
l <- ggplot(seals, aes(long, lat))
```
- l + geom\_contour(aes(z = z))**  
x, y, z, alpha, colour, group, linetype, size, weight

#### continuous bivariate distribution

- ```
h <- ggplot(diamonds, aes(carat, price))
```

- h + geom\_bin2d(binwidth = c(0.25, 500))**  
x, y, alpha, color, fill, linetype, size, weight

- h + geom\_density2d()**  
x, y, alpha, colour, group, linetype, size

- h + geom\_hex()**  
x, y, alpha, colour, fill, size

#### continuous function

- ```
i <- ggplot(economics, aes(date, unemploy))
```

- i + geom\_area()**  
x, y, alpha, color, fill, linetype, size

- i + geom\_line()**  
x, y, alpha, color, group, linetype, size

- i + geom\_step(direction = "hv")**  
x, y, alpha, color, group, linetype, size

#### visualizing error

- ```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

- j + geom\_crossbar(fatten = 2)**  
x, y, ymax, ymin, alpha, color, fill, group, linetype, size

- j + geom\_errorbar()**, x, ymax, ymin, alpha, color, group, linetype, size, width (also `geom_errorbarh()`)

- j + geom\_linerange()**  
x, ymin, ymax, alpha, color, group, linetype, size

- j + geom\_pointrange()**  
x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

#### maps

- ```
data <- data.frame(murder = USAArrests$Murder,
state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))
```

- k + geom\_map(aes(map\_id = state), map = map)**  
+ expand\_limits(x = map\$long, y = map\$lat), map\_id, alpha, color, fill, linetype, size



# Ice cream analysis



The data were collected on 200 high school students and are scores on various tests, including a video game and a puzzle. The outcome measure in this analysis is the student's favorite flavor of ice cream – vanilla, chocolate or strawberry- from which we are going to see what relationships exists with video game scores (**video**), puzzle scores (**puzzle**) and gender (**female**).

| Variable name    | Variable                                                  | Data type        |
|------------------|-----------------------------------------------------------|------------------|
| <b>Id</b>        | Identity of the student                                   | Nominal          |
| <b>female</b>    | Gender (0: Male, 1:Female)                                | Binary           |
| <b>ice_cream</b> | Favorite Flavor (1: Vanilla, 2: Chocolate, 3: Strawberry) | Nominal          |
| <b>video</b>     | Score on the video game                                   | Scale/Continuous |
| <b>puzzle</b>    | Score on the puzzle                                       | Scale/Continuous |