

(RE)INTRODUCING SOME STATISTICAL TERMS

feat. teacup giraffes



Rachel Warnock with help from
tinystats.github.io/teacups-giraffes-and-statistics

21.02.2022

GENERAL CONCEPTS

- ▶ Intro to the Normal Distribution
- ▶ Mean, median, mode
- ▶ Spread of the data
- ▶ Standard Error



TYPES OF DATA

CONTINUOUS

measured data, can have ∞ values within possible range.



I AM 3.1" TALL

I WEIGH 34.16 grams

DISCRETE

OBSERVATIONS CAN ONLY EXIST
AT LIMITED VALUES, OFTEN COUNTS.



I HAVE 8 LEGS
and
4 SPOTS!

NOMINAL

UNORDERED DESCRIPTIONS



ORDINAL

ORDERED DESCRIPTIONS



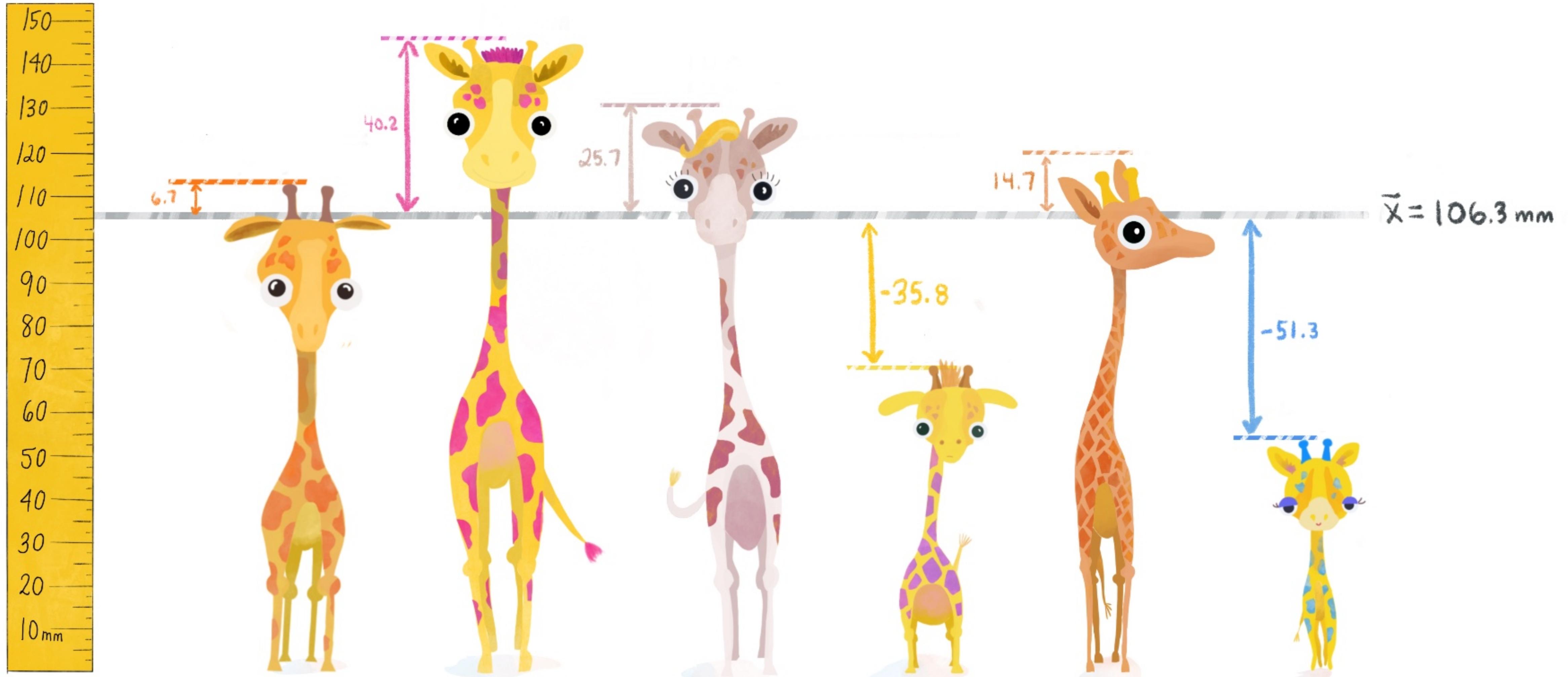
BINARY

ONLY 2 MUTUALLY EXCLUSIVE OUTCOMES



WHAT TYPES OF DATA ARE USED IN THE TESTS IN YOUR PAPER?

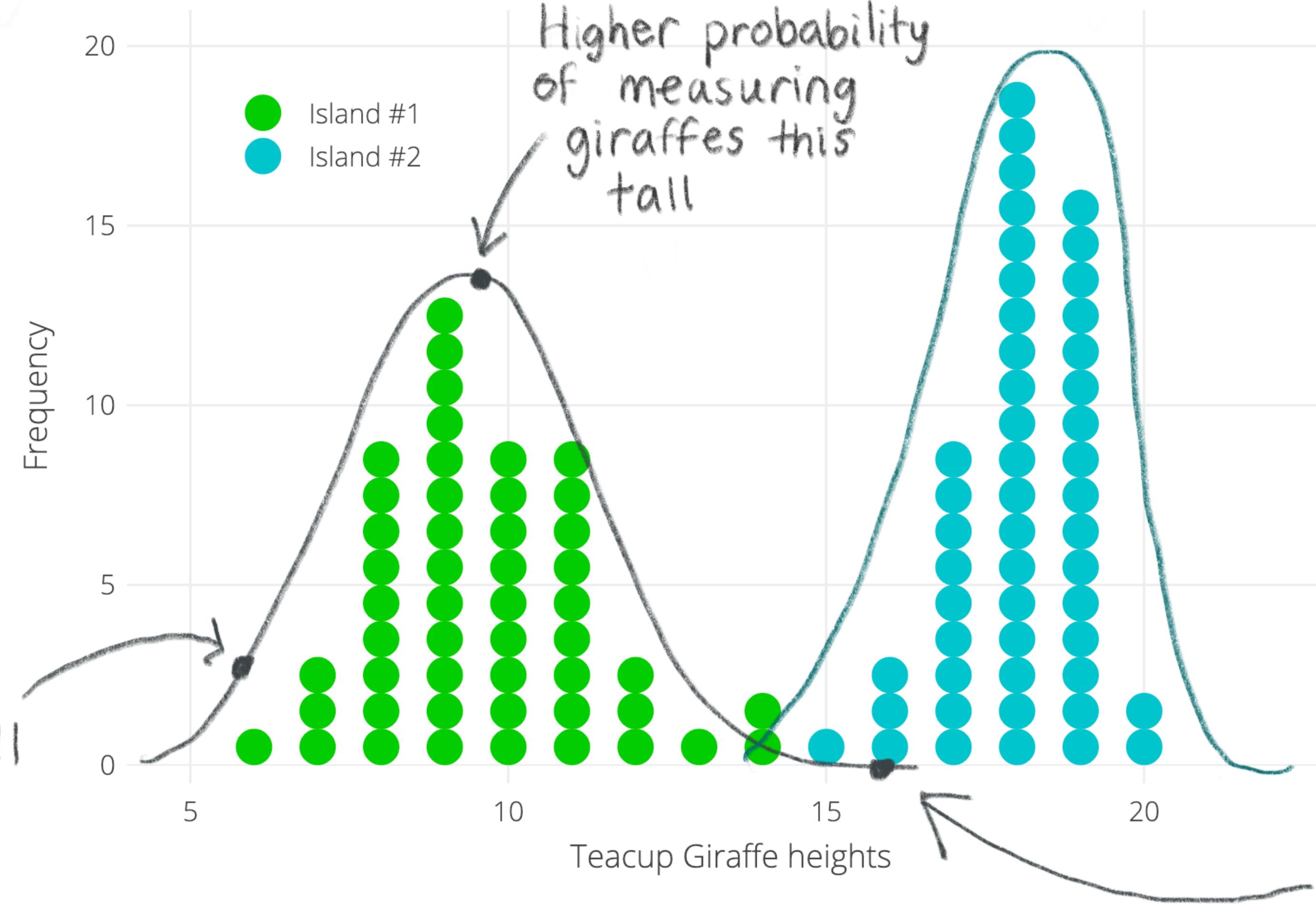
IMAGINE WE'VE VISITED TWO ISLANDS OF TINY GIRAFFES





[teacups-giraffes-and-statistics](#)

Encountering a giraffe this small would be rare



the normal distribution

SAMPLE VS POPULATION

A **sample**, e.g. giraffes from each island (or different coloured birds), is a subset of a **population**.

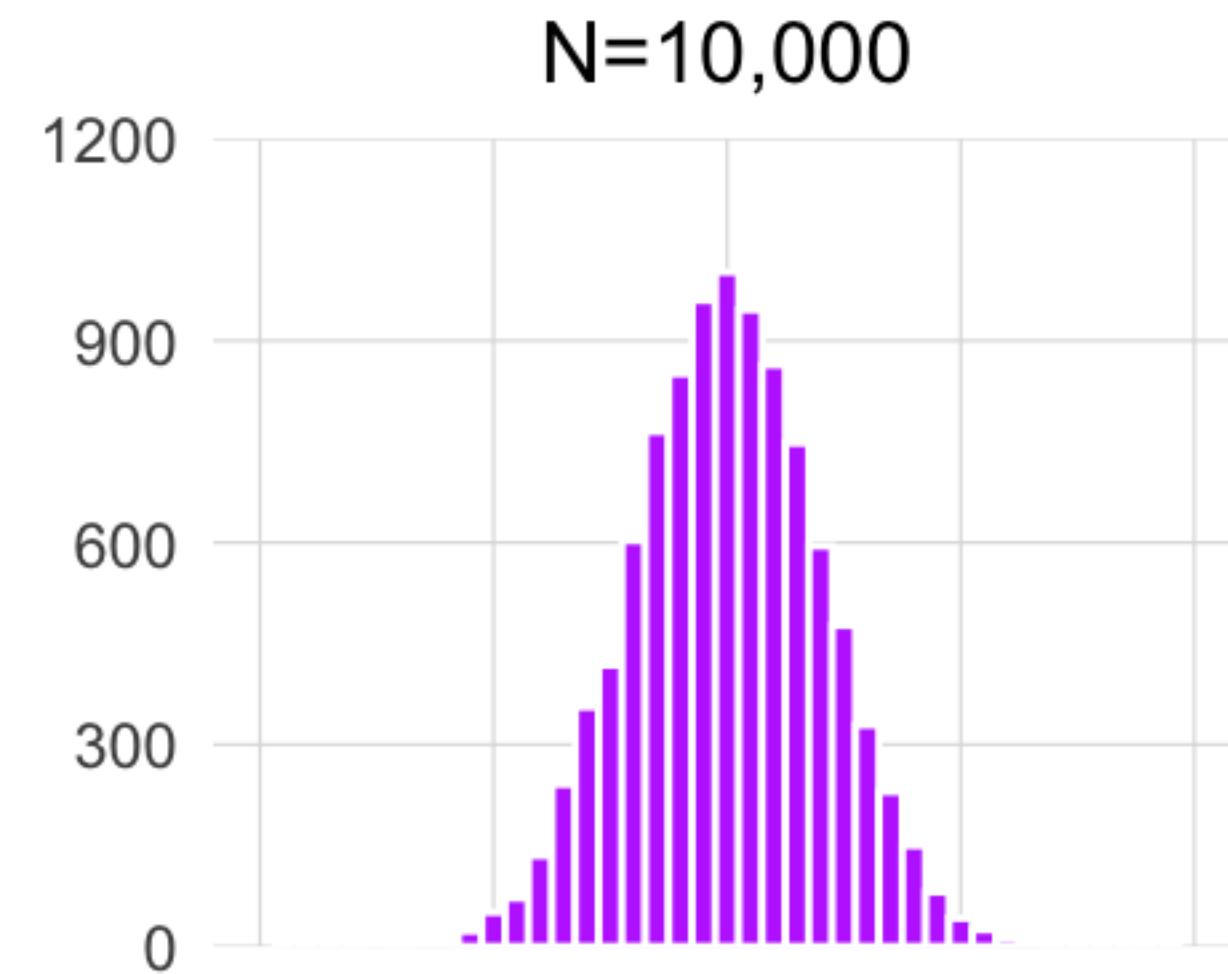
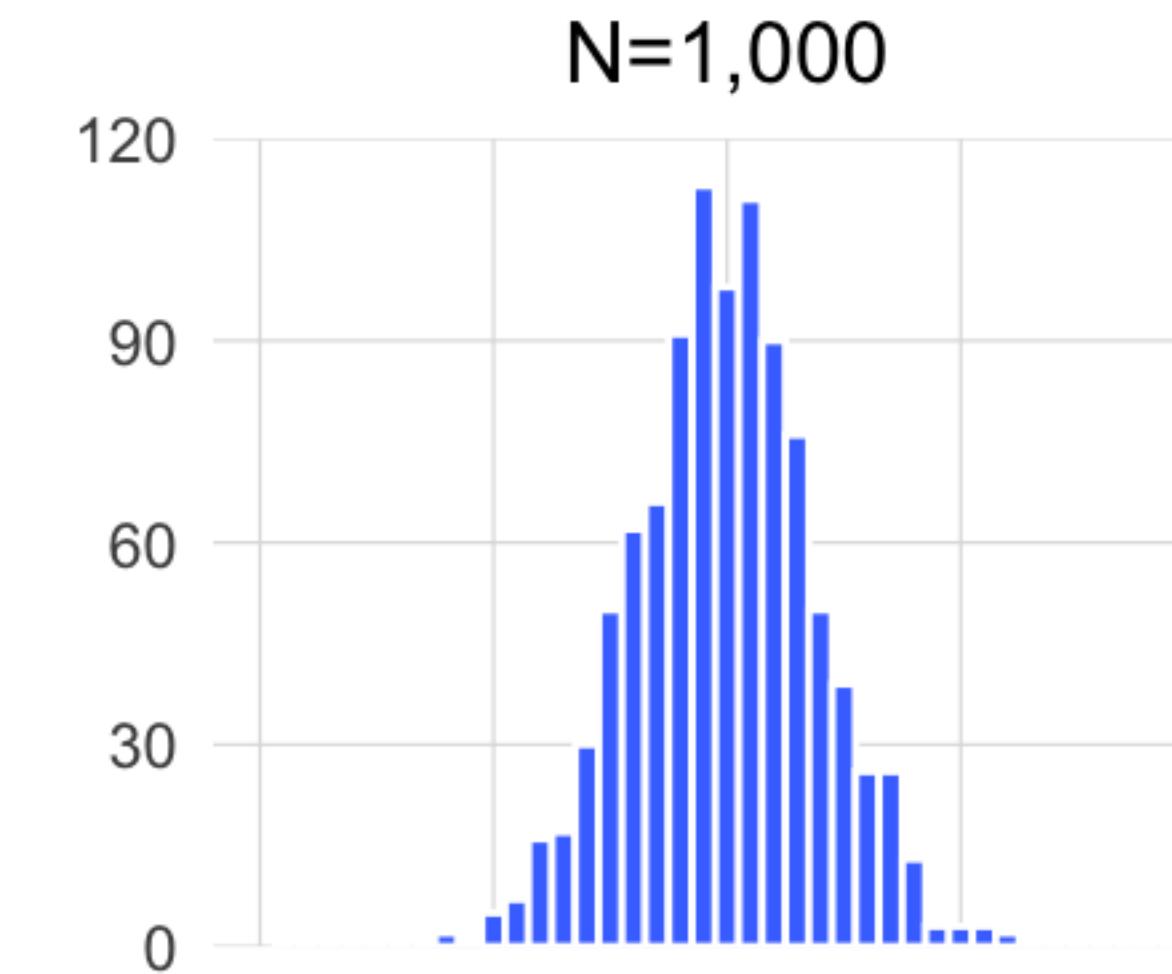
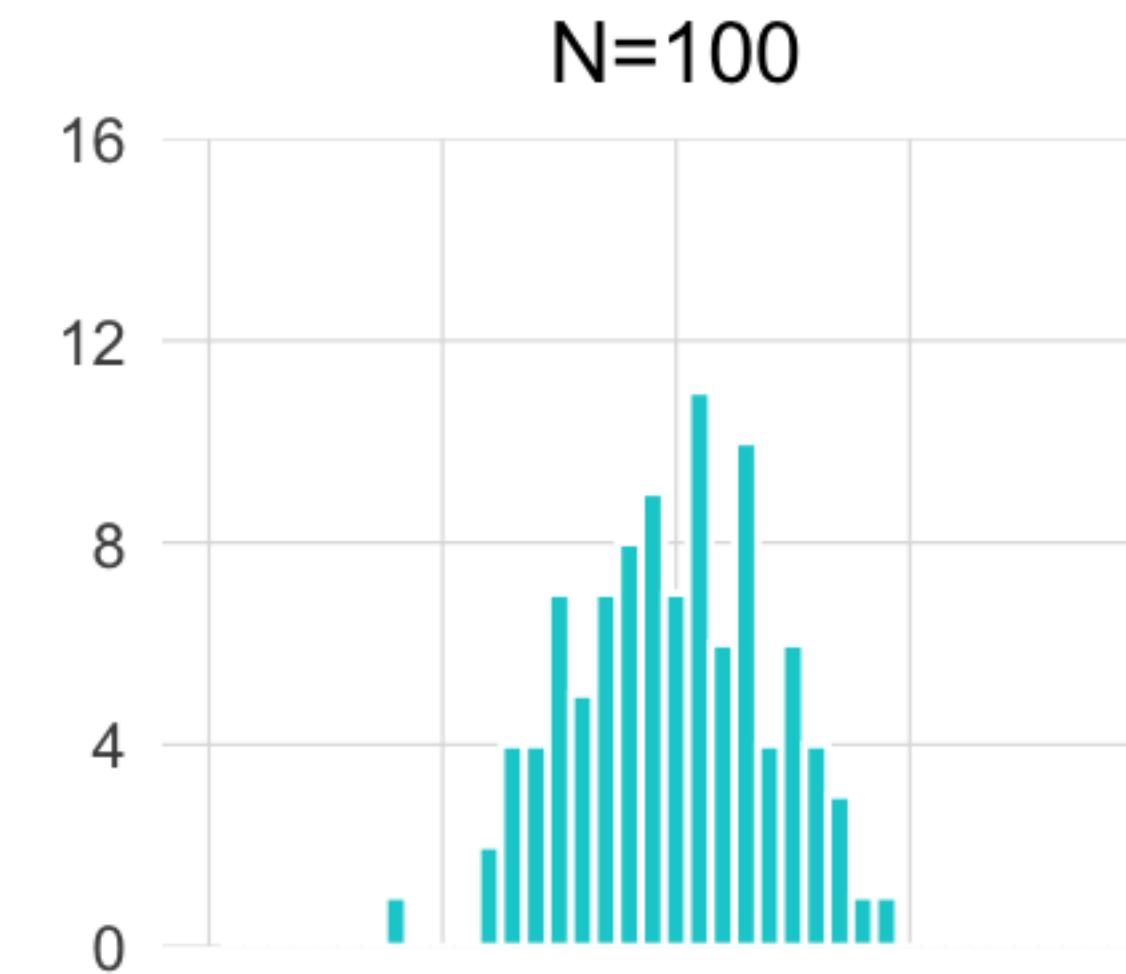
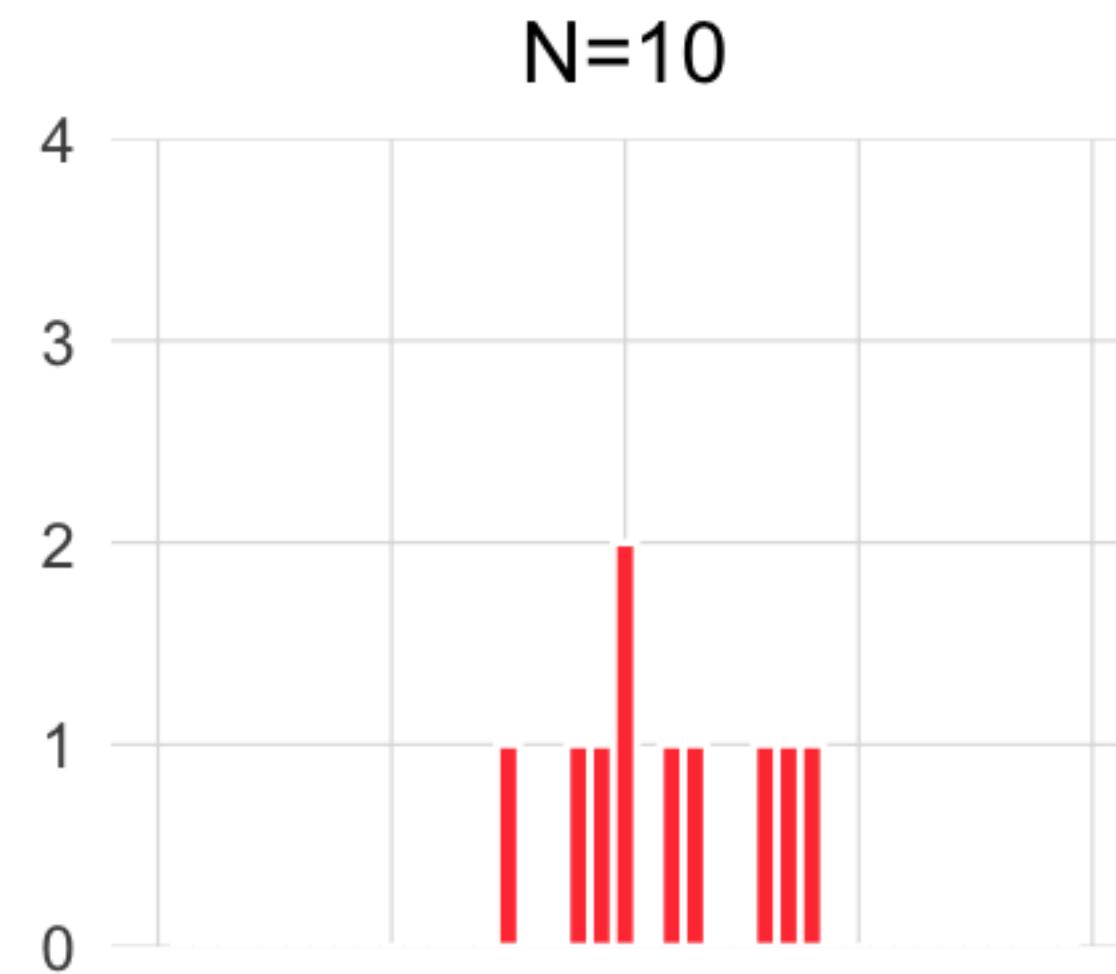
The population is defined as all available observations. In this case, all existing giraffes on one of the islands or all existing coloured birds.

Different symbols are used to represent the mean for each of these.

SAMPLE VS. POPULATION

| | | |
|-------------|---------------------------------|----------------------------------|
| \bar{X} | Sample mean | Yes calculated from the raw data |
| μ | True population mean | Almost never known for sure |
| $\hat{\mu}$ | Estimate of the population mean | Yes identical to the sample mean |

SAMPLE SIZE



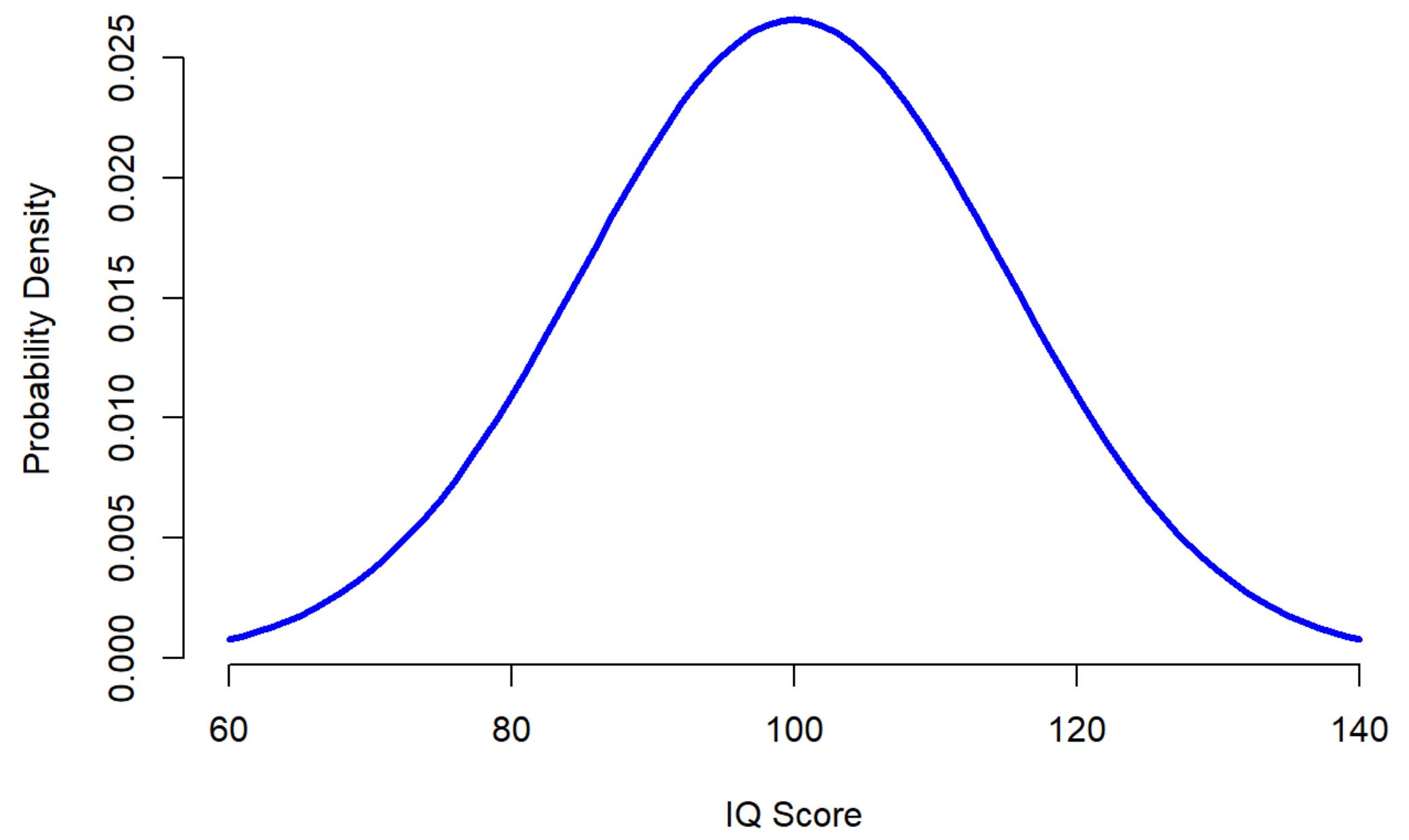
the larger our sample is, the more of the population it will include

SAMPLING

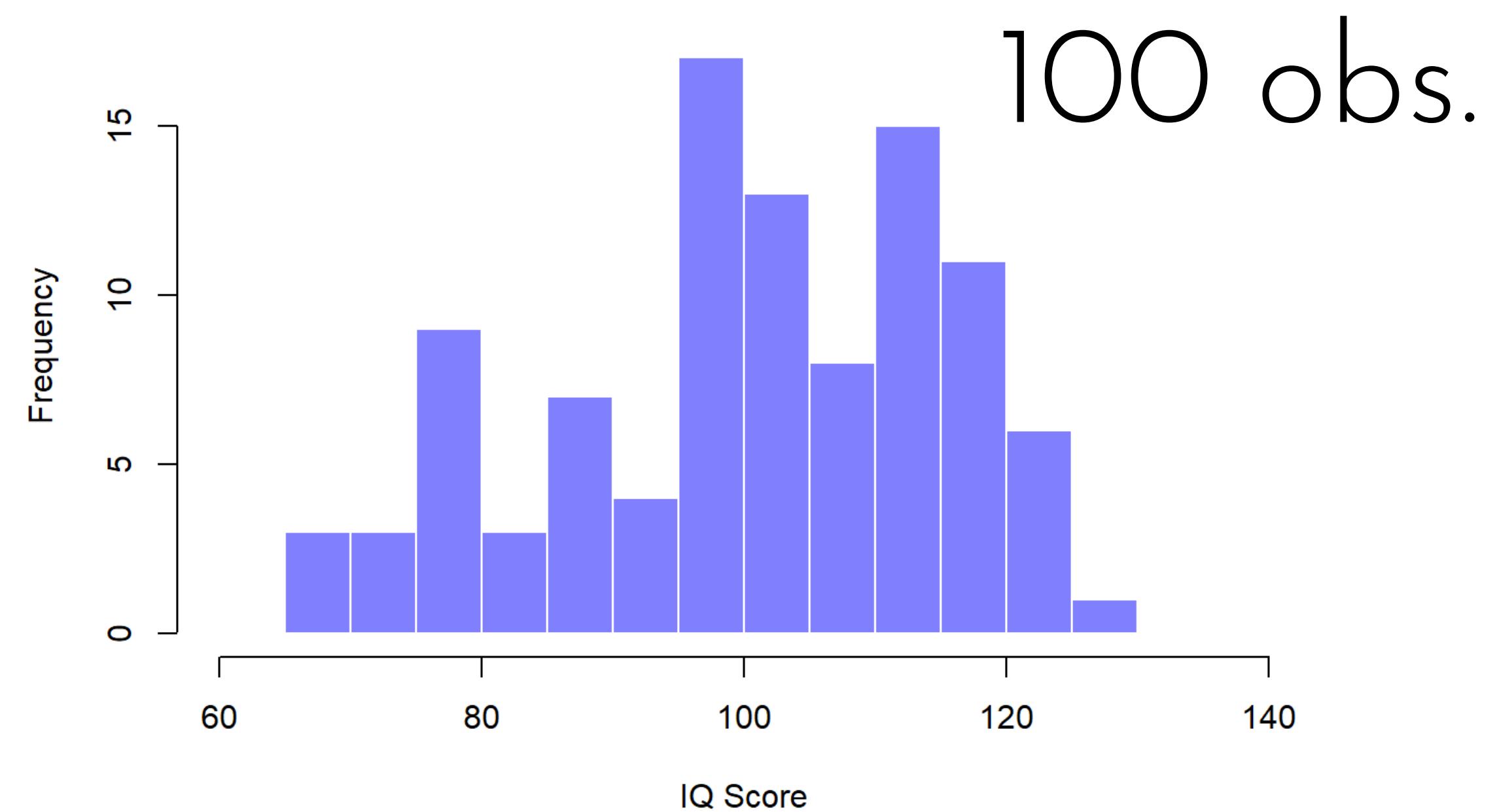
A procedure in which every member of the population has the same chance of being selected is called a **simple random sample**.

In reality, most samples are not random. Many samples are biased.

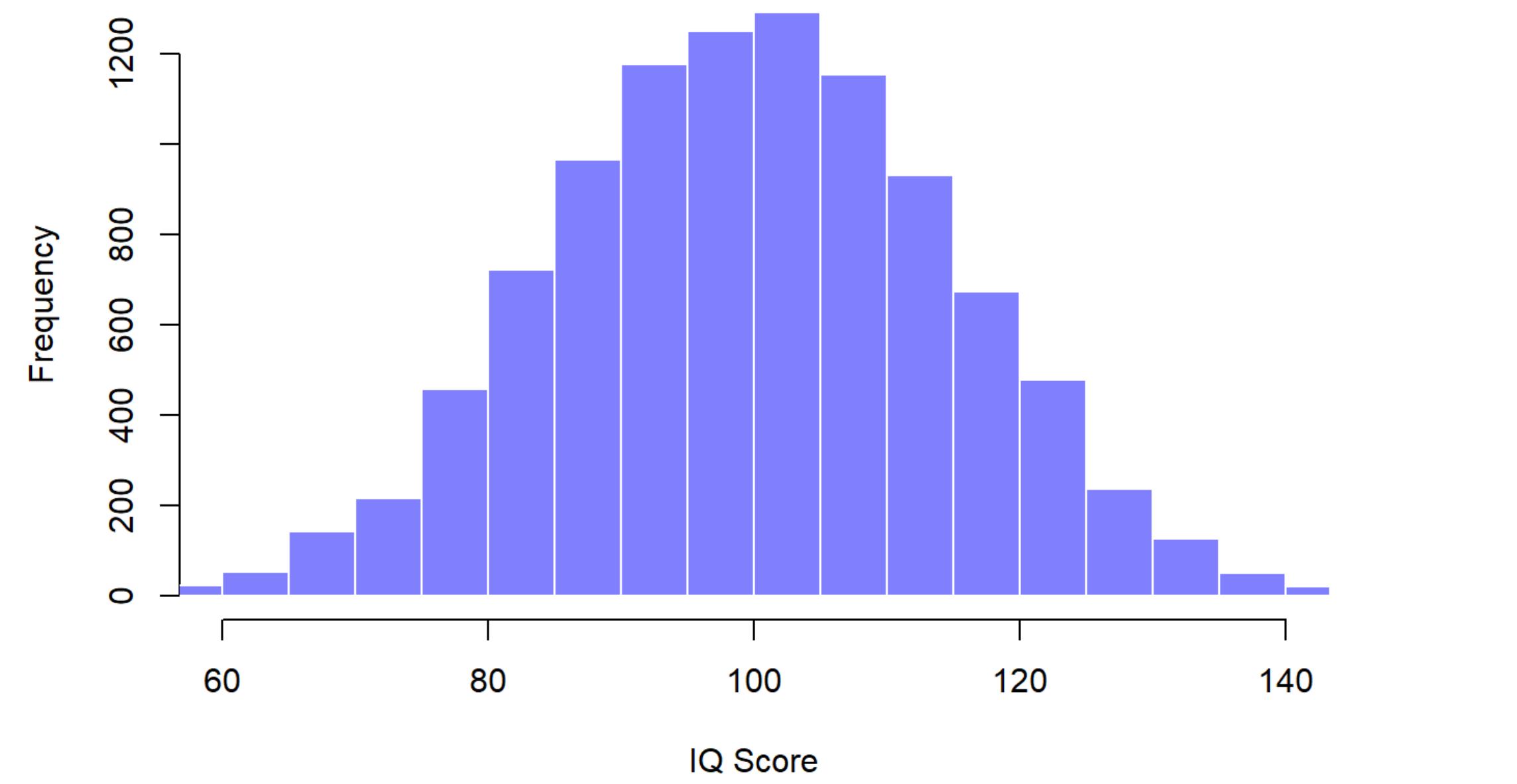
POPULATION PARAMETERS



Learning statistics with R



100 obs.

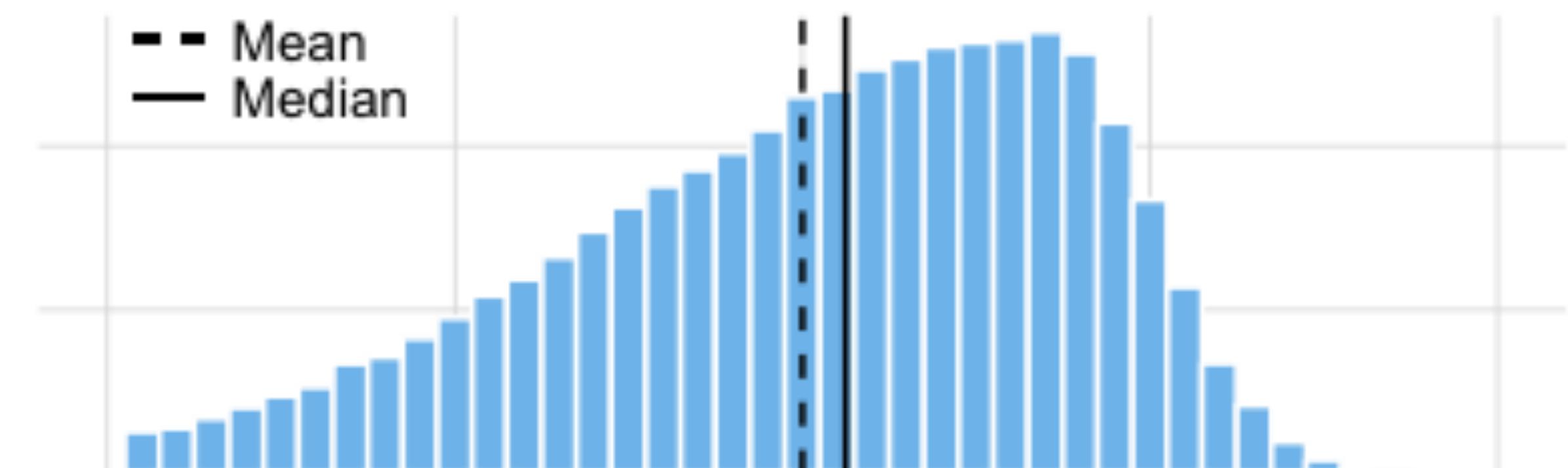


10000 obs.

MEASURES OF CENTRALITY

Mean – the average and the measure of centrality.

Median – the value in the middle of the data set.

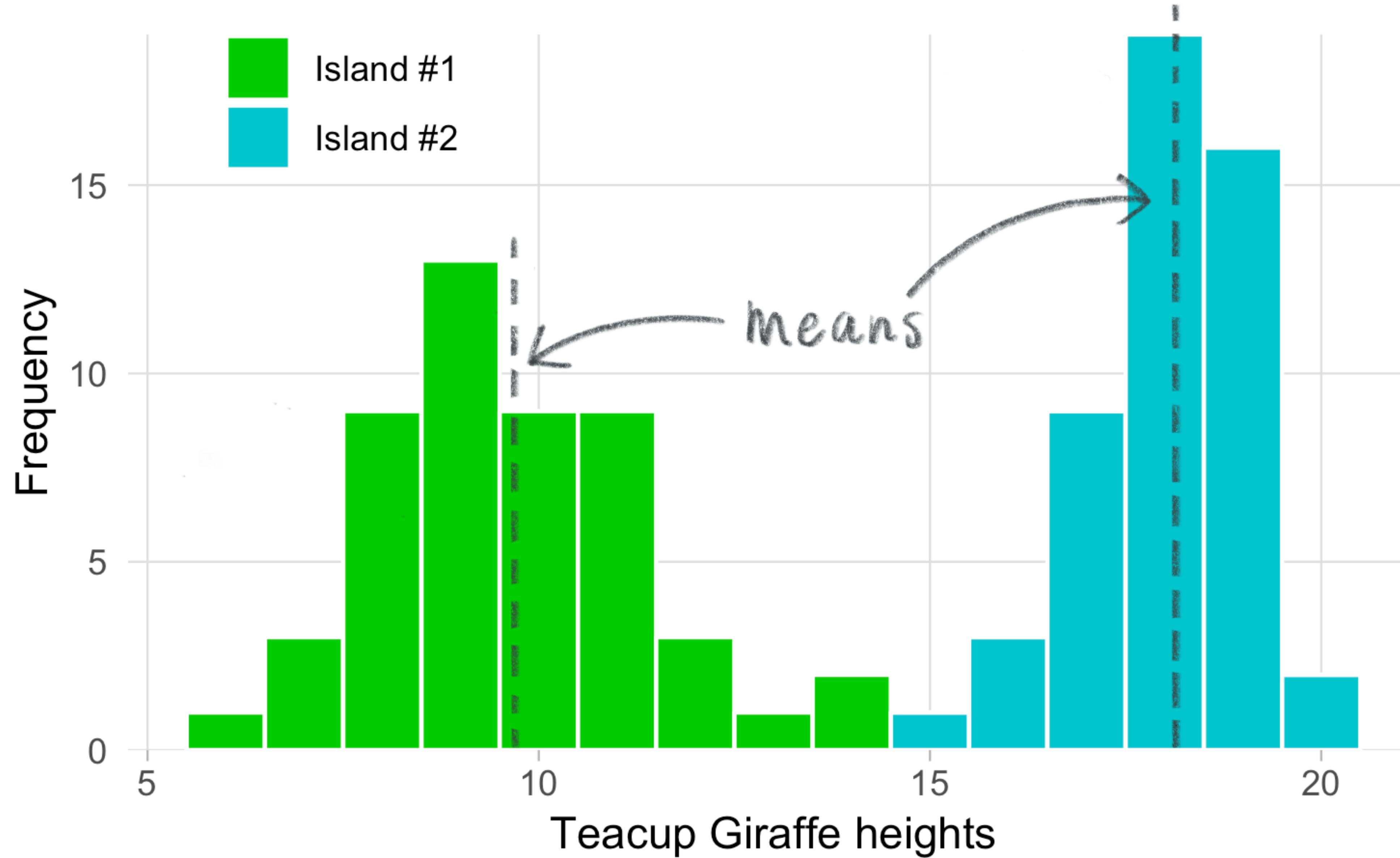


MEASURES OF CENTRALITY

Mean – the average and the measure of centrality.

Median – the value in the middle of the data set.

Mode – the value (e.g. height of giraffes, bird colour) that occurs most frequently.



$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

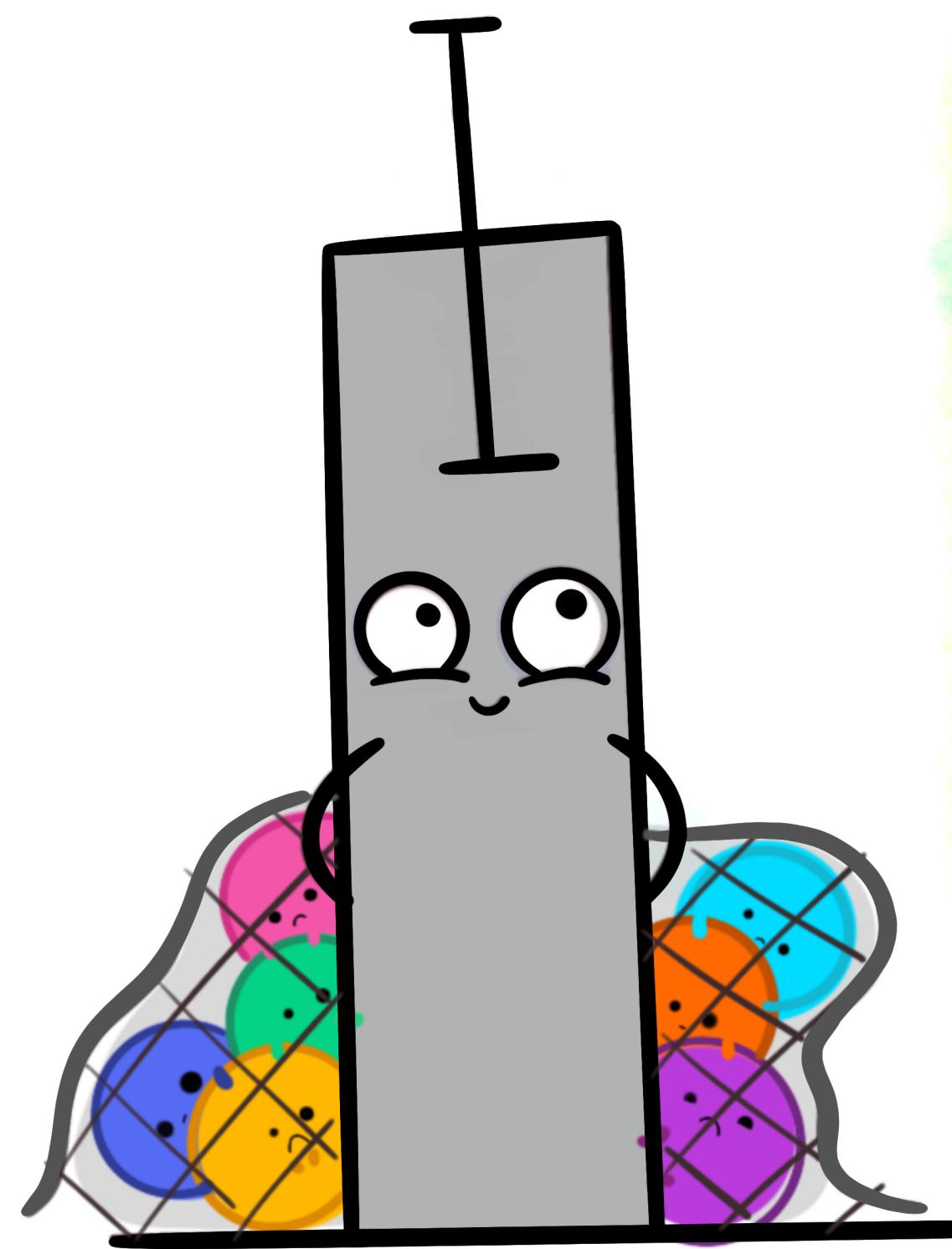
"add up" →

start with the first term ↗

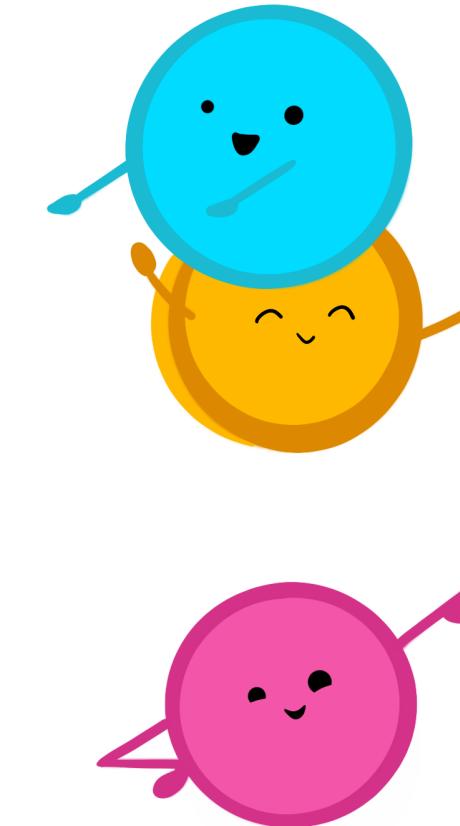
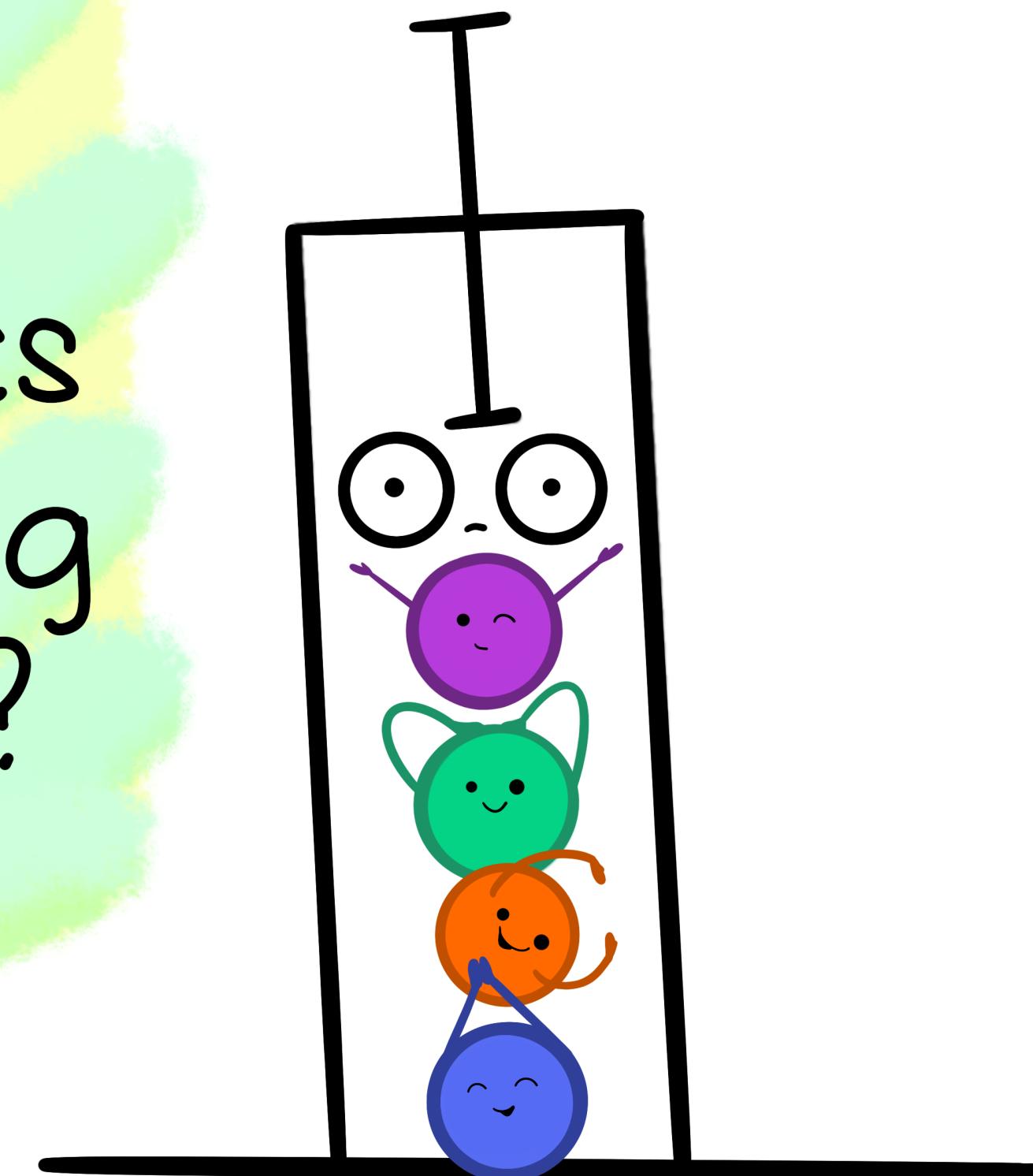
keep adding until you reach the n^{th} (last) term

x_i ← each term (the heights)

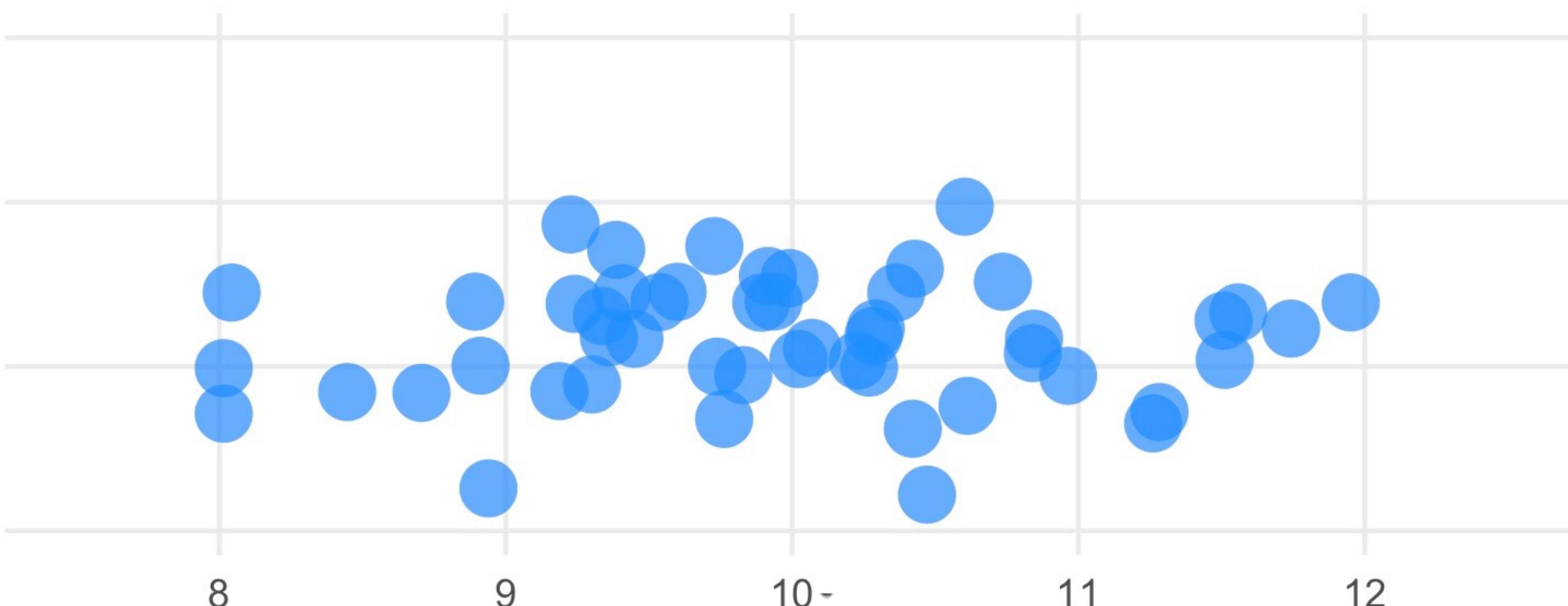
n ↗ total sample size



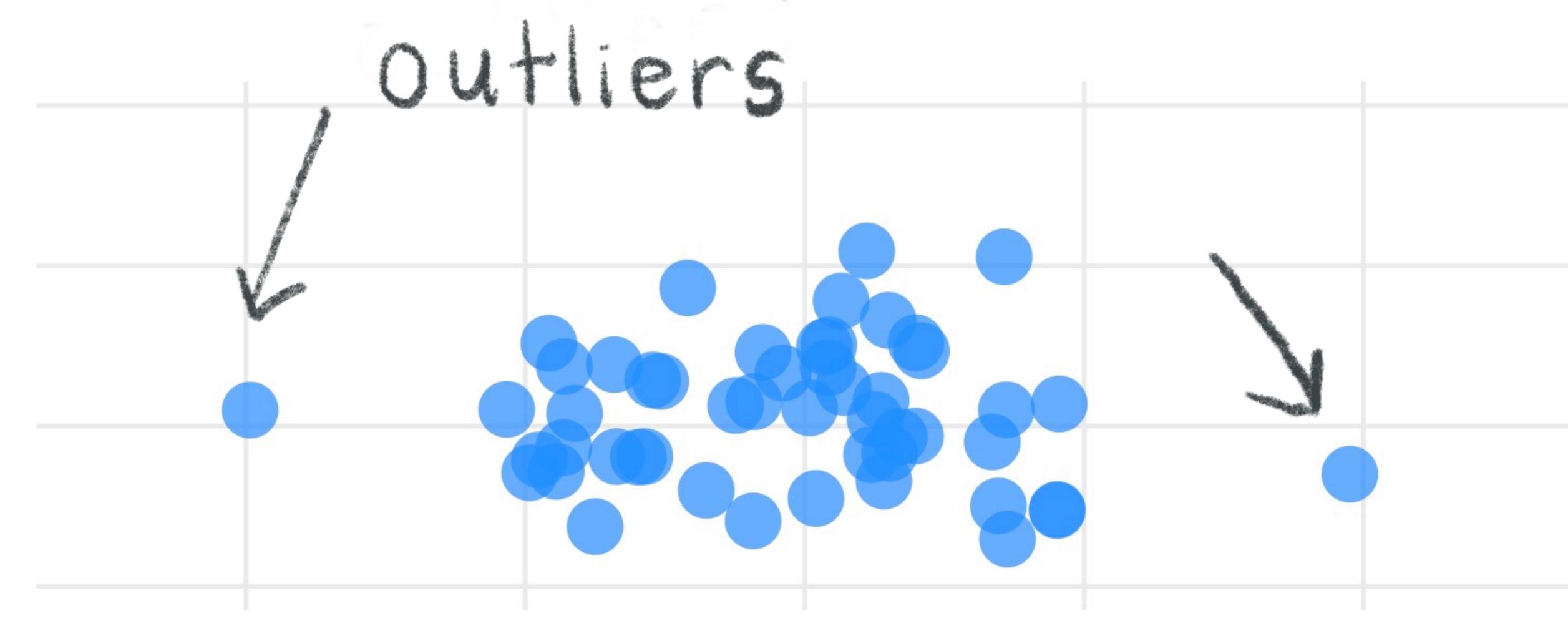
are your
summary statistics
hiding something
interesting?



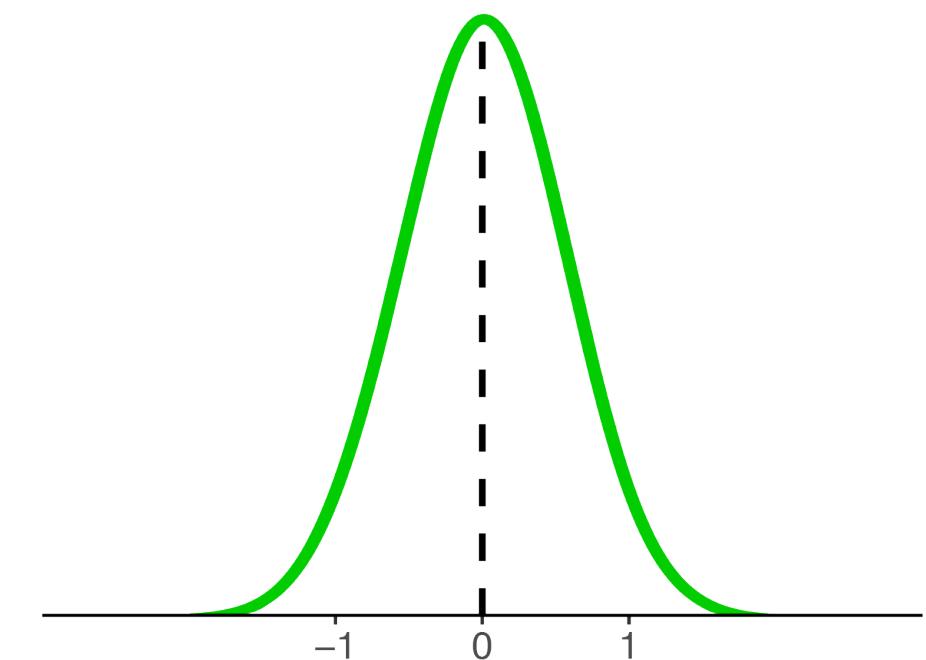
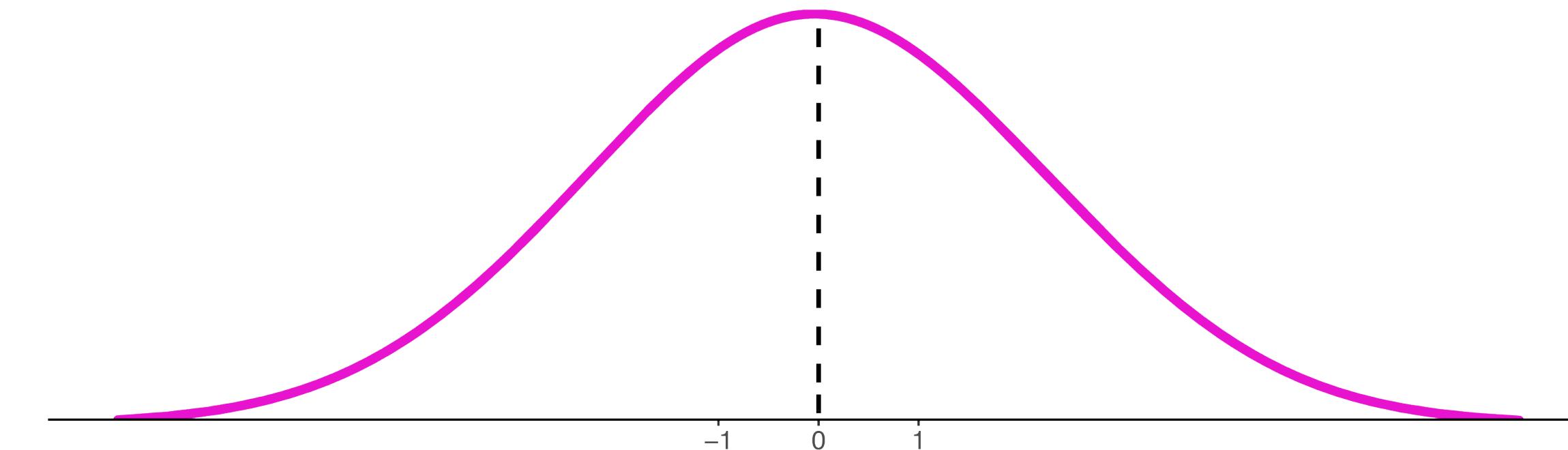
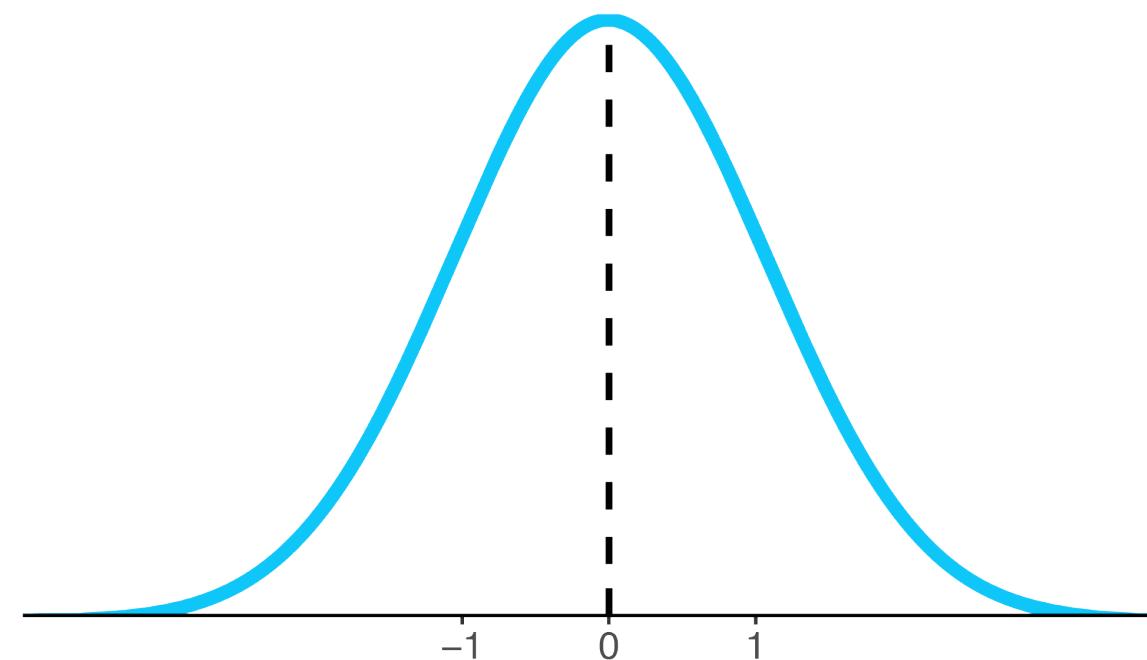
SPREAD OF THE DATA - THE RANGE



but range = 4
for both



SPREAD OF THE DATA - VARIANCE

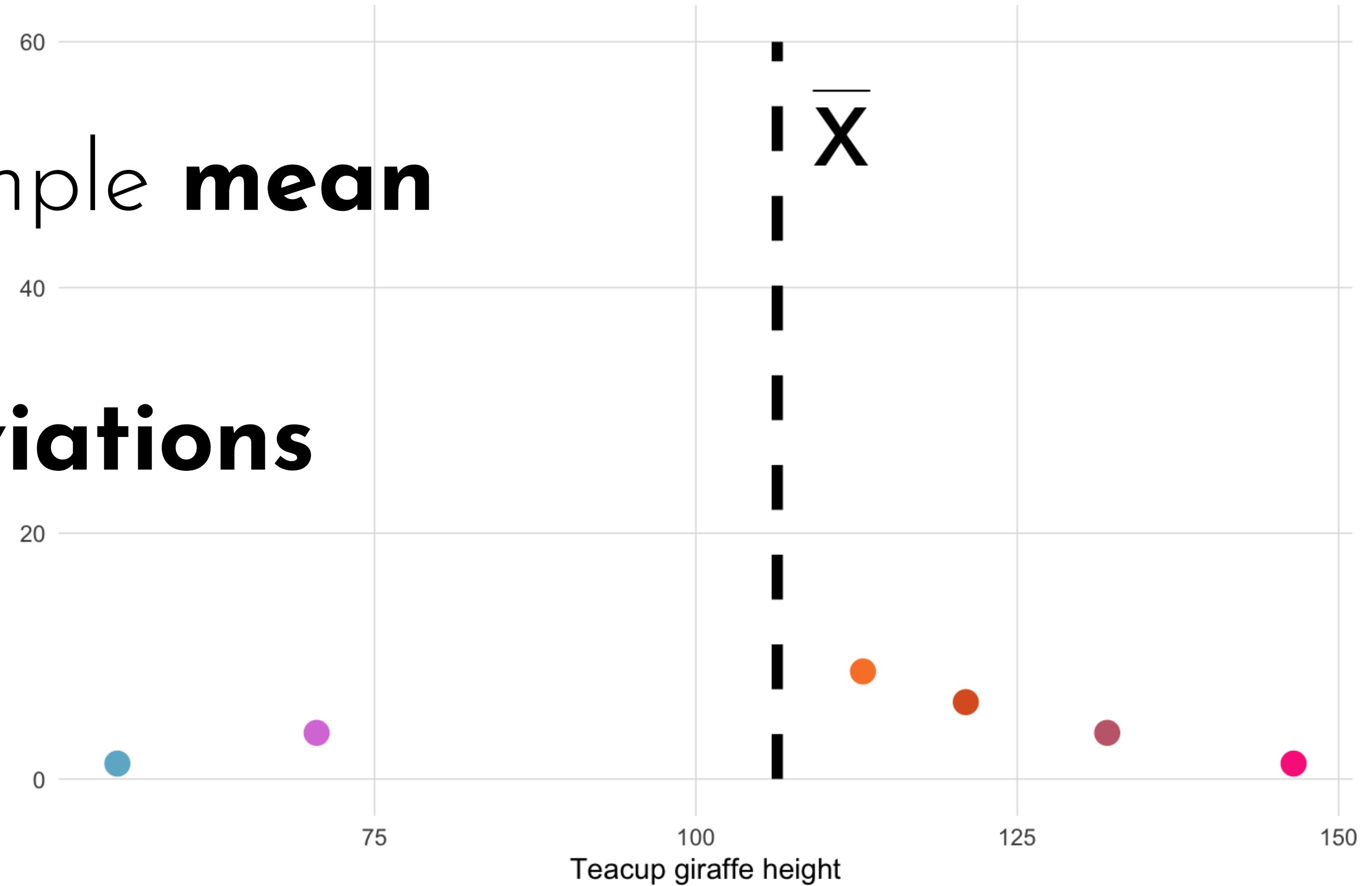


HOW DO WE CALCULATE THIS?

CALCULATING VARIANCE

1. Calculate the sample **mean**

2. **Square the deviations**
from the mean



CALCULATING VARIANCE

3. Calculate the sum of squares

A diagram illustrating the calculation of the sum of squares. On the left, six colored squares (blue, pink, purple, light red, orange, and a small orange square) are shown with plus signs between them, followed by an equals sign. To the right of the equals sign is a large red square with a double equals sign at its bottom edge. This visualizes the mathematical operation of squaring each value and then summing them up.

$$\sum_{i=1}^N (x_i - \mu)^2$$

CALCULATING VARIANCE

4. Calculate the average of the squared differences from the mean.

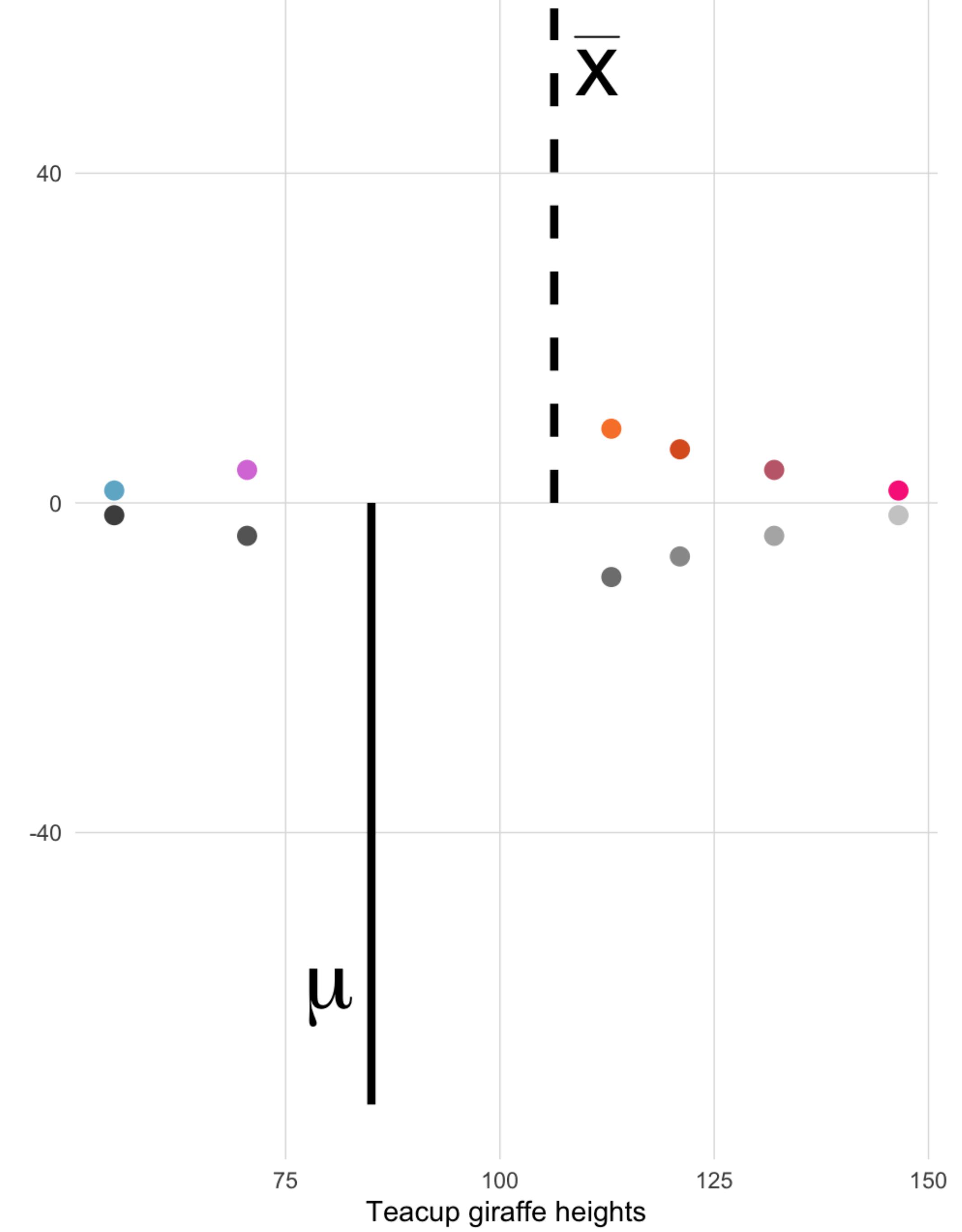
$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

STANDARD DEVIATION

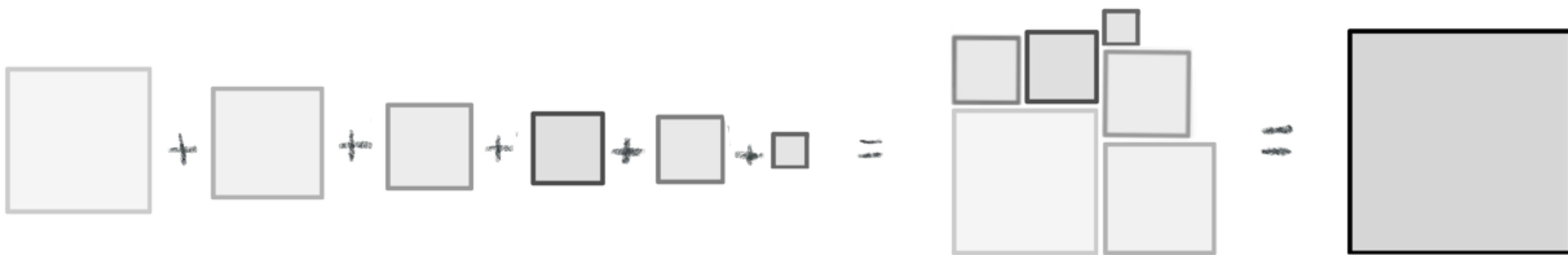
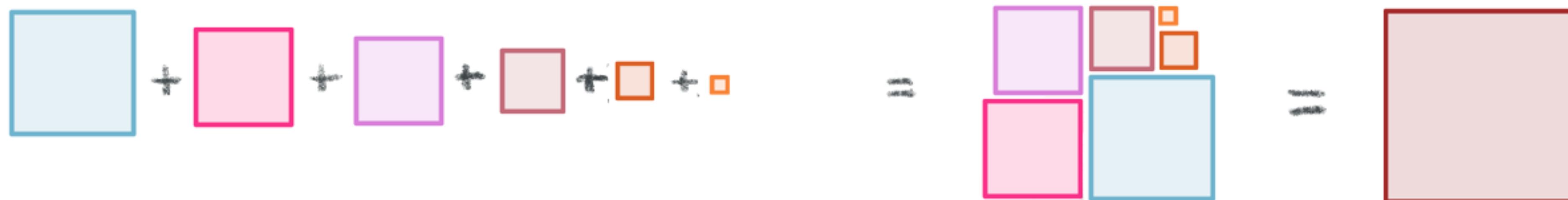
Variance is not easily interpretable → “unsquare”
the variance to return to the data’s original units.

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

POPULATION VERSUS SAMPLE



POPULATION VERSUS SAMPLE



SOLUTION : N - 1

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

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

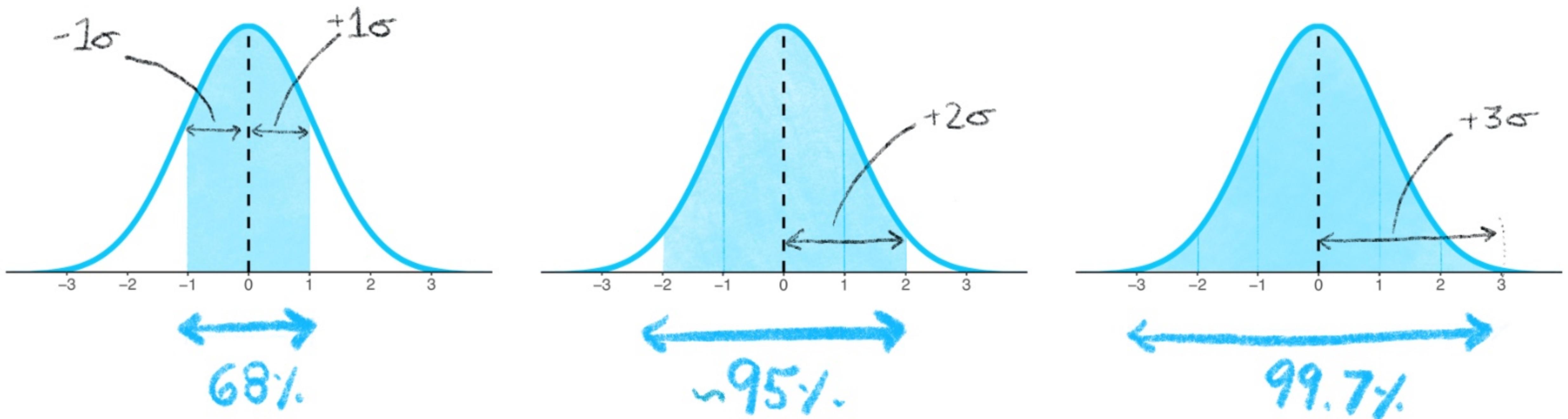
$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

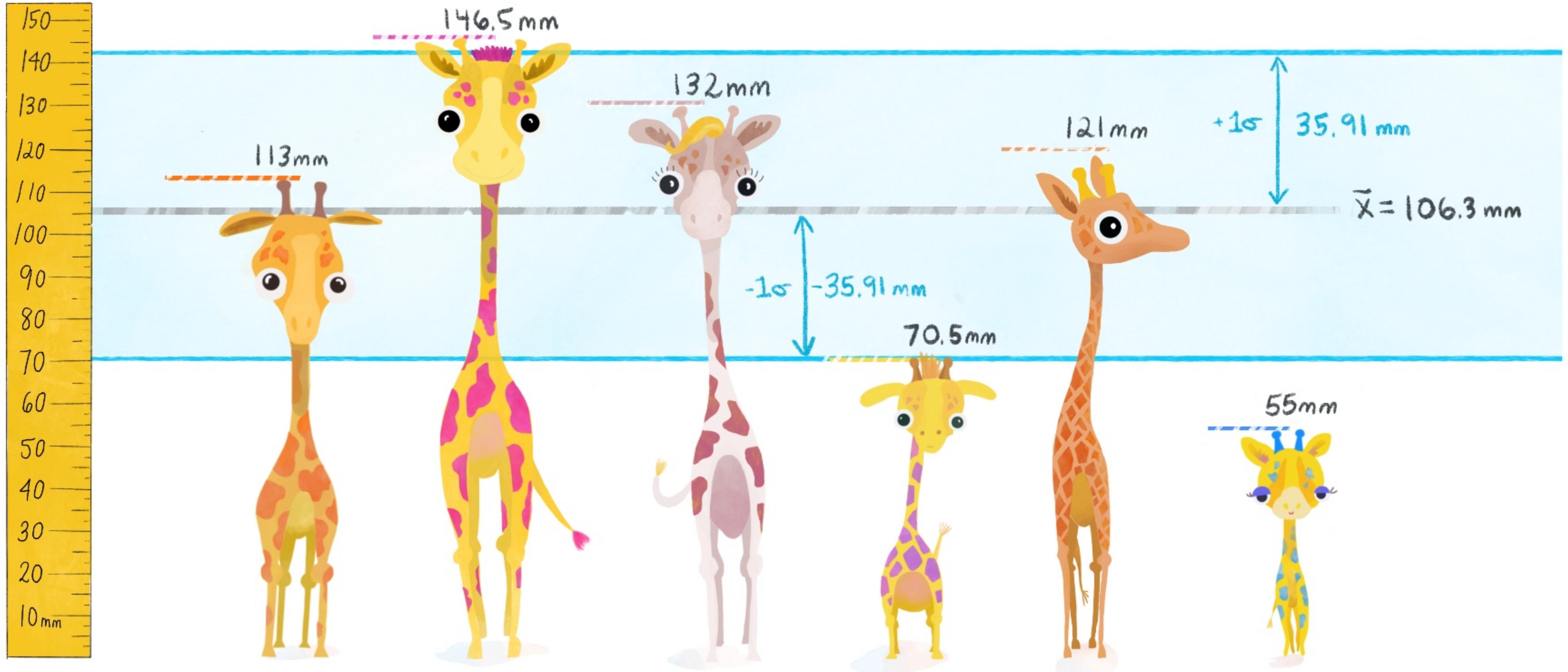
Also note the change in symbols

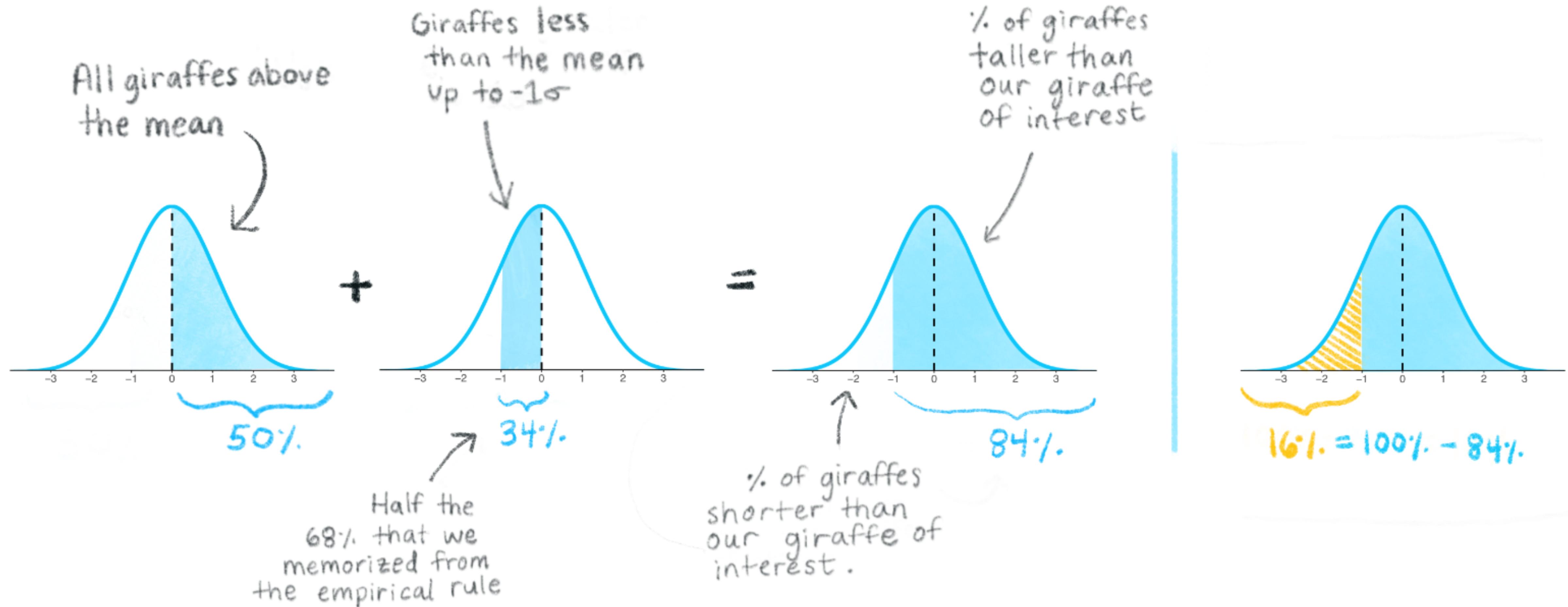
STEPS SUMMARY

1. Calculate the sample **mean**
2. Square the **deviations from the mean**
3. Calculate the **sum of squares**
4. Calculate **average squared differences**
5. Apply the N-1 correction
6. Compare your results

INTERPRETING THE STANDARD DEVIATION

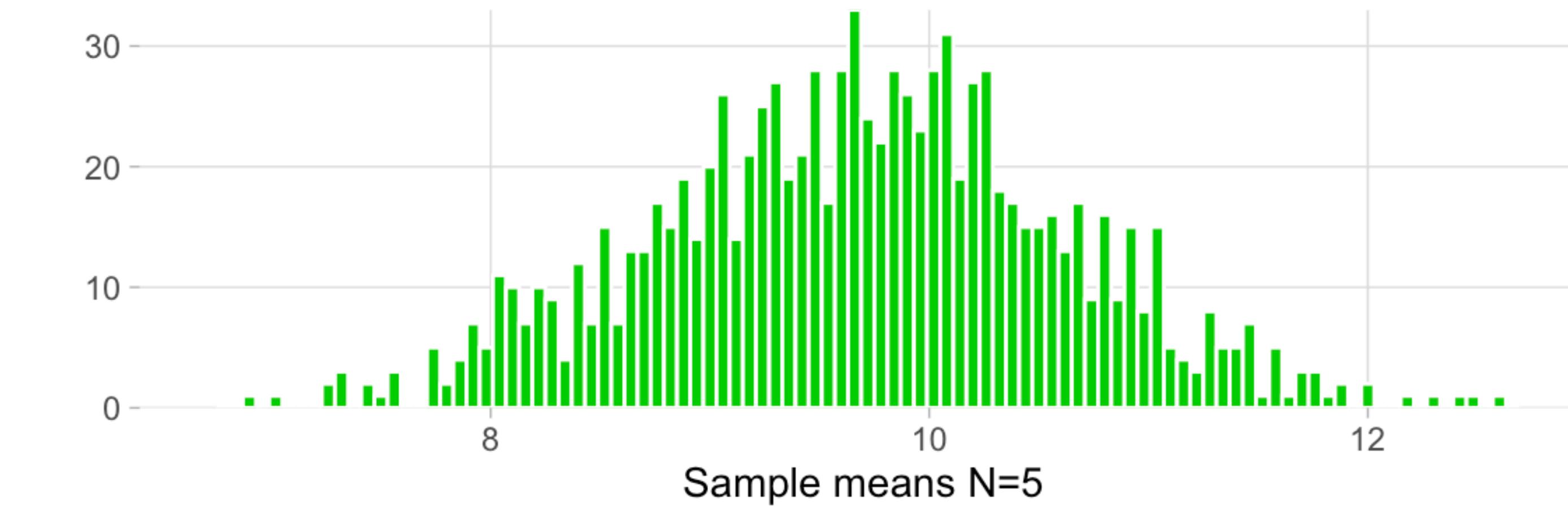
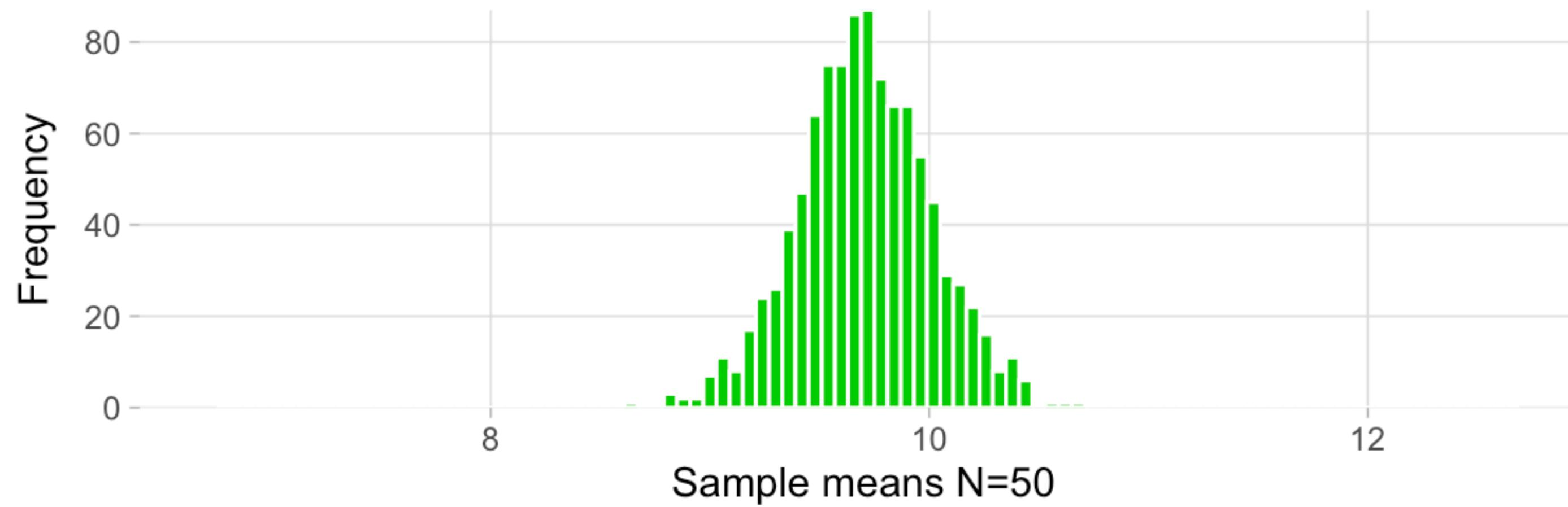
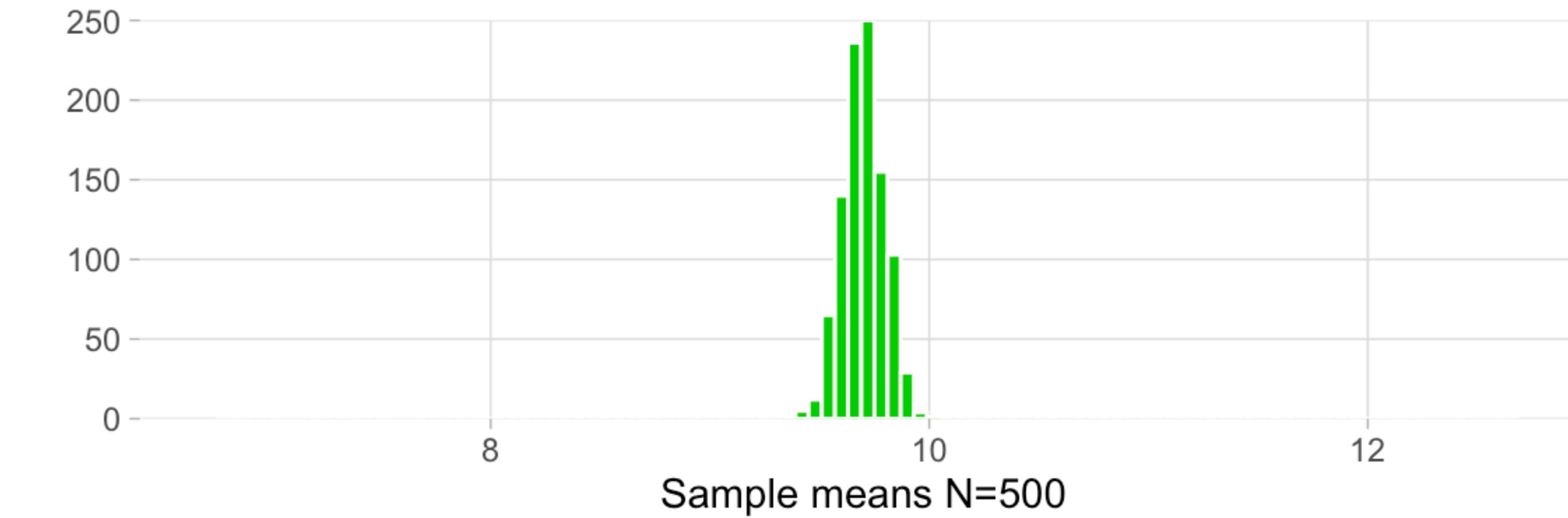




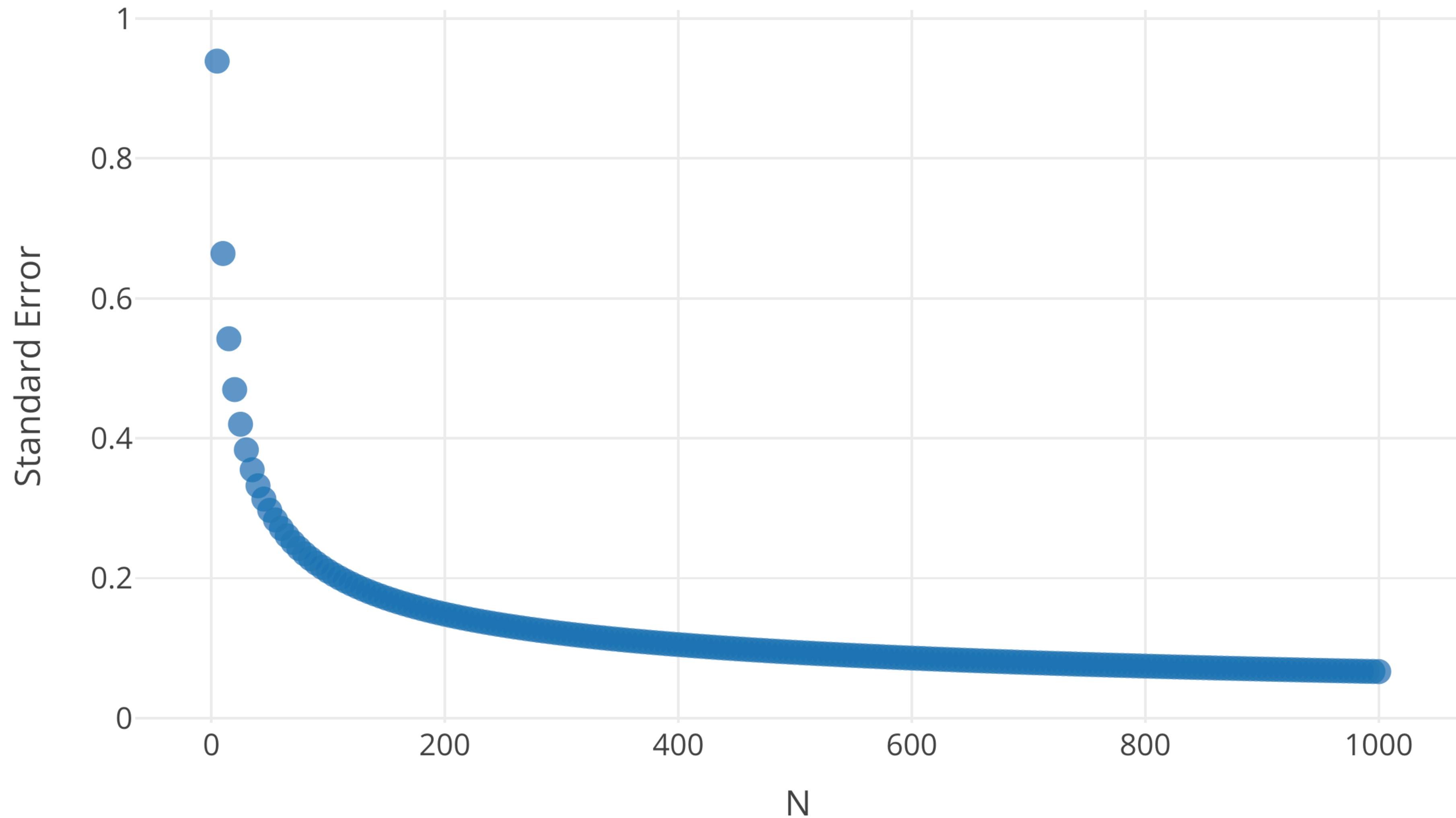


SAMPLE SIZE

distributions of
the mean based
on random
samples of
different sizes



STANDARD ERROR



STANDARD ERROR

$$\mathbf{SE} = \frac{\sigma}{\sqrt{n}}$$

- a measure of the statistical accuracy of an estimate
- equal to the standard deviation of the theoretical distribution of a large population of equivalent estimates

FURTHER READING

tinystats.github.io/teacups-giraffes-and-statistics



QUESTIONS?