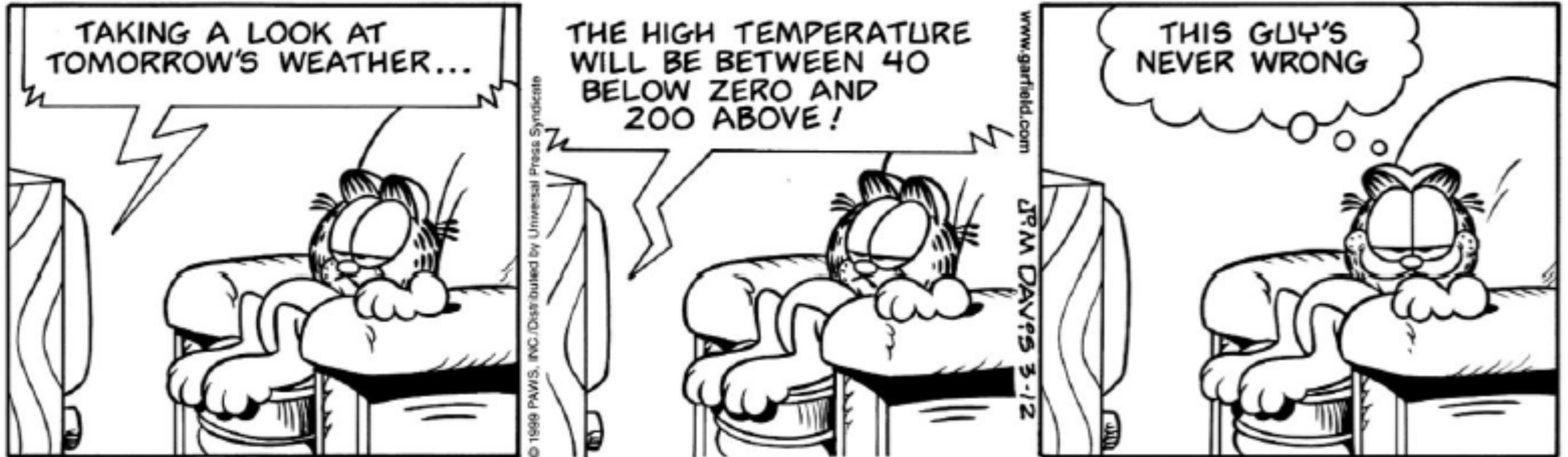


# Lecture 06

10.12.21



# Refresher Quiz

**The heights of men in a certain population follow a normal distribution with mean 69.7 inches and standard deviation 2.8 inches.**

**a) If a man is chosen at random from the population, find the probability that he will be more than 72 inches tall.**

**b) If two men were chosen at random from the population, find the probability that (i) both of them will be more than 72 inches tall; (ii) their mean height will be more than 72 inches.**

The heights of men in a certain population follow a normal distribution with mean 69.7 inches and standard deviation 2.8 inches.

a) If a man is chosen at random from the population, find the probability that he will be more than 72 inches tall.

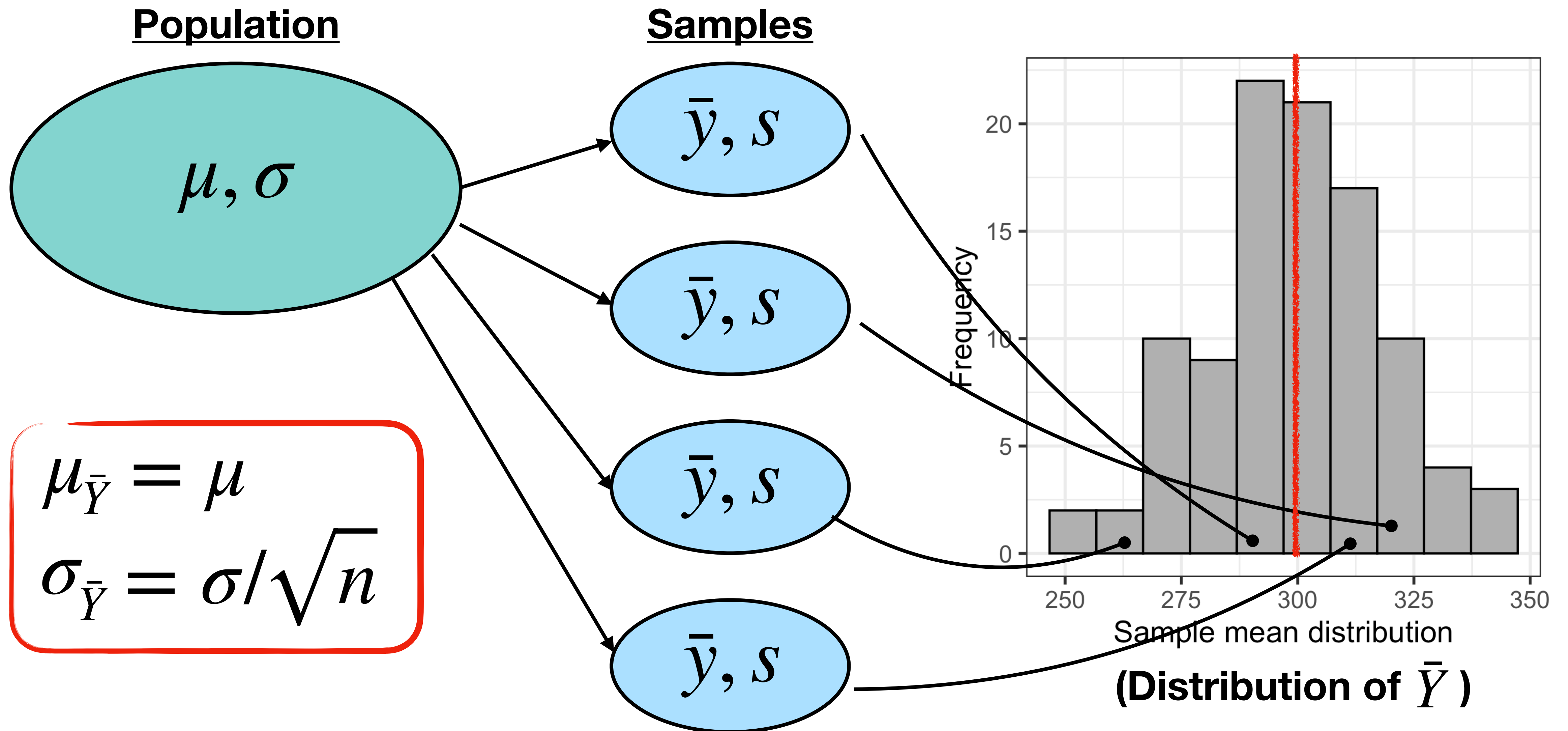
```
> pnorm(72, 69.7, 2.8, lower.tail = F) = 20.5%
```

b) If two men were chosen at random from the population, find the probability that (i) both of them will be more than 72 inches tall; (ii) their mean height will be more than 72 inches.

$P(\text{both} > 72 \text{ in tall}) = P(> 72) * P(> 72) = 0.205 * 0.205 = 4.20\%$

```
> pnorm(72, 69.7, 2.8/sqrt(2), lower.tail = F) = 12.26%
```

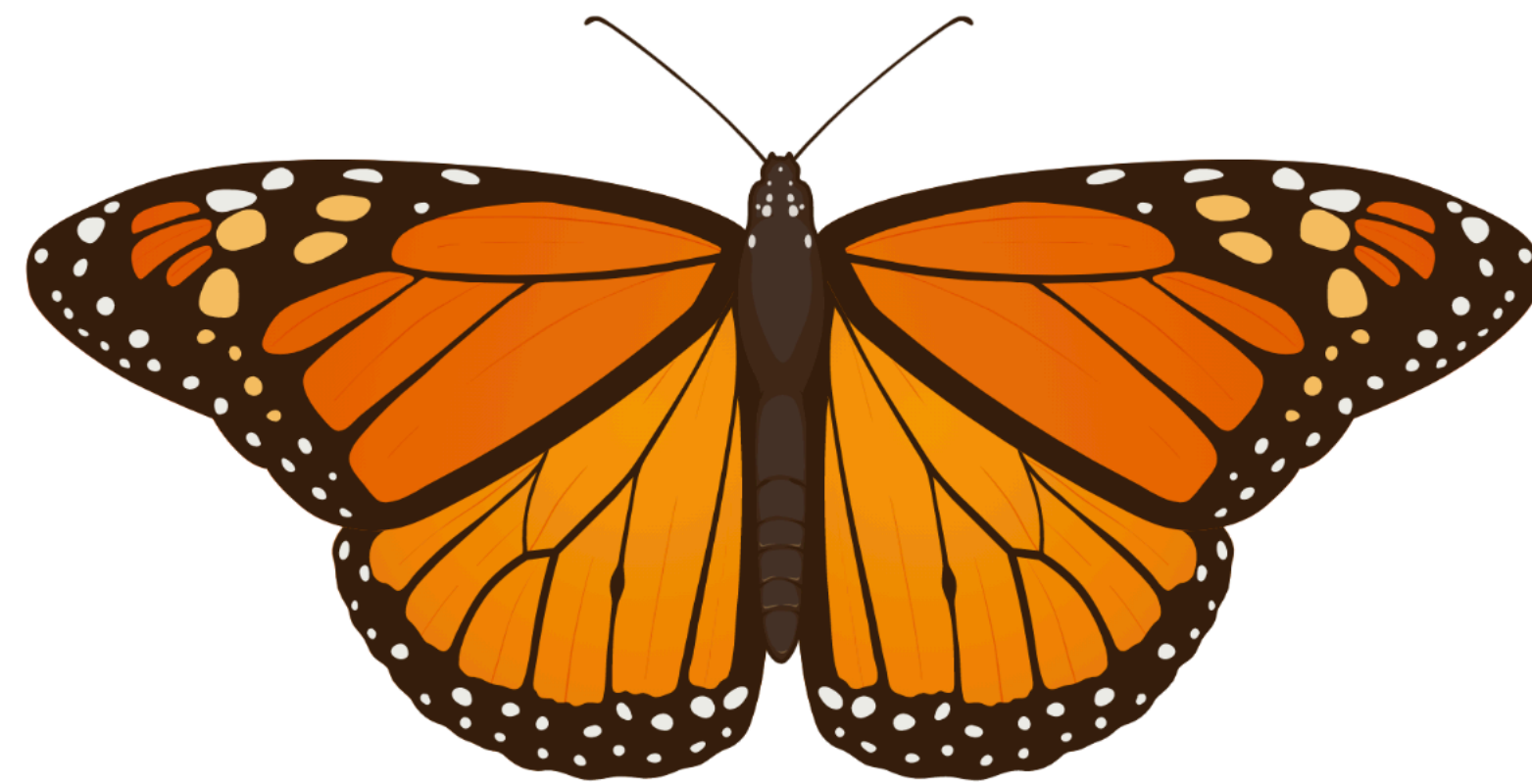
# Summary: sampling distribution



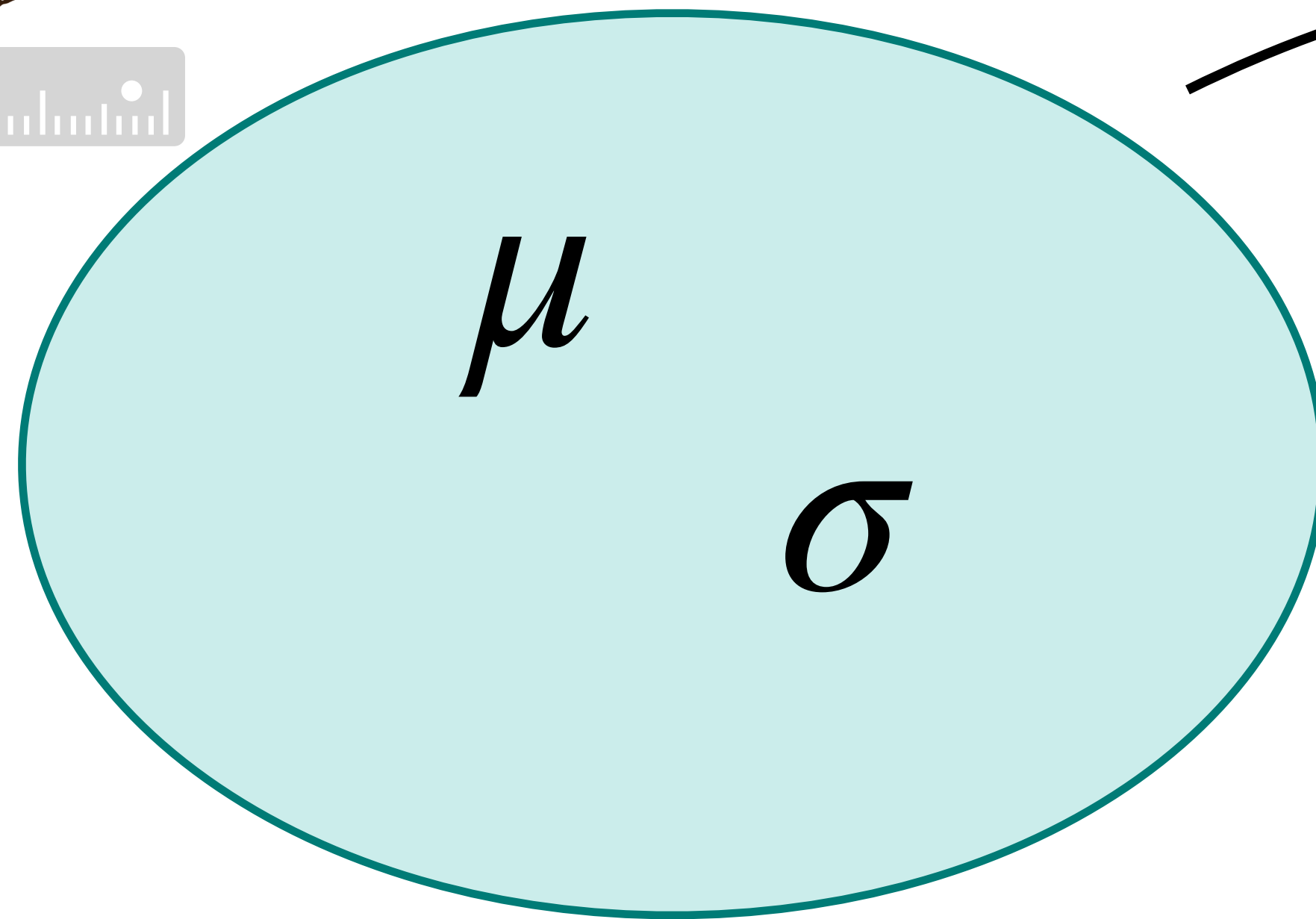
# Statistical estimation

- We view our data as a **random sample from a population** and use the information about our data to *infer* facts about the population
- Goals:
  - (1) Determine estimate of some feature of the population (i.e. mean)
  - (2) Assess the precision of the estimate

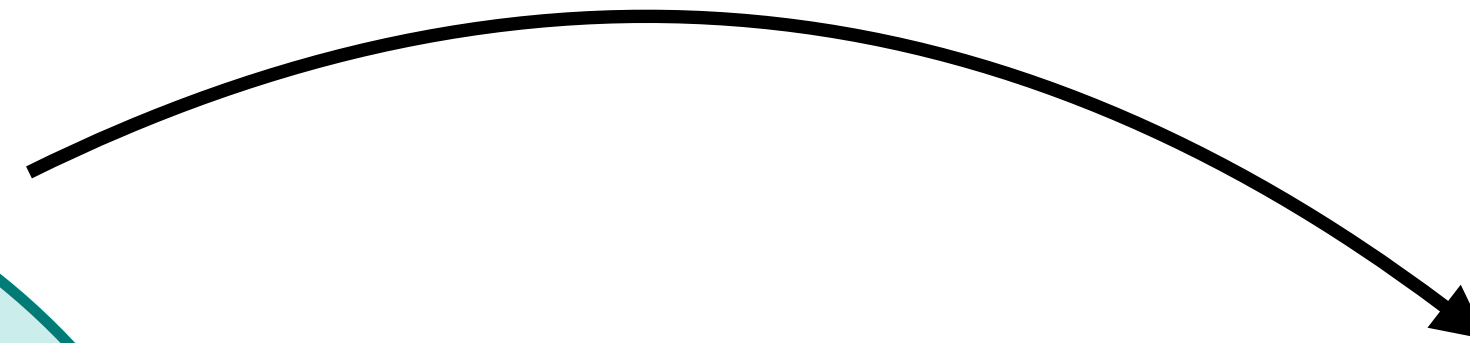
# Statistical estimation



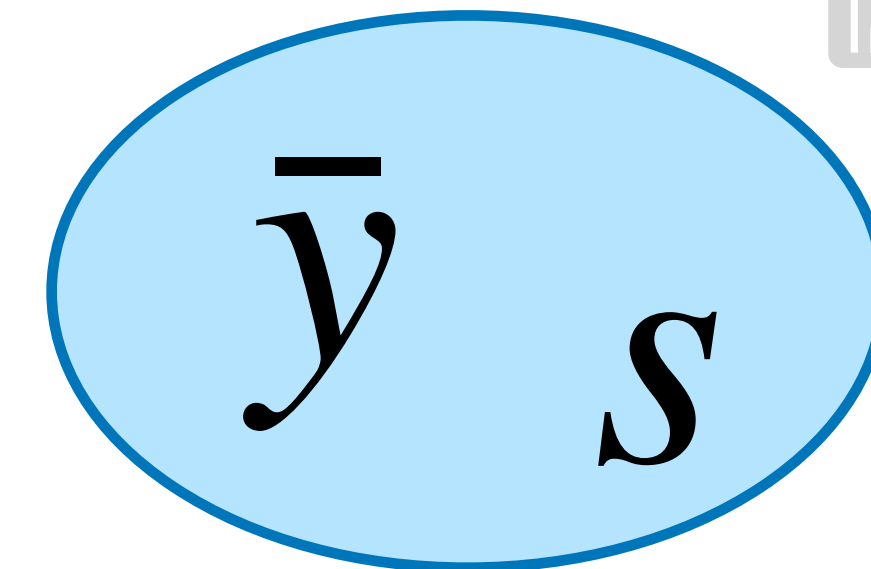
# Statistical estimation



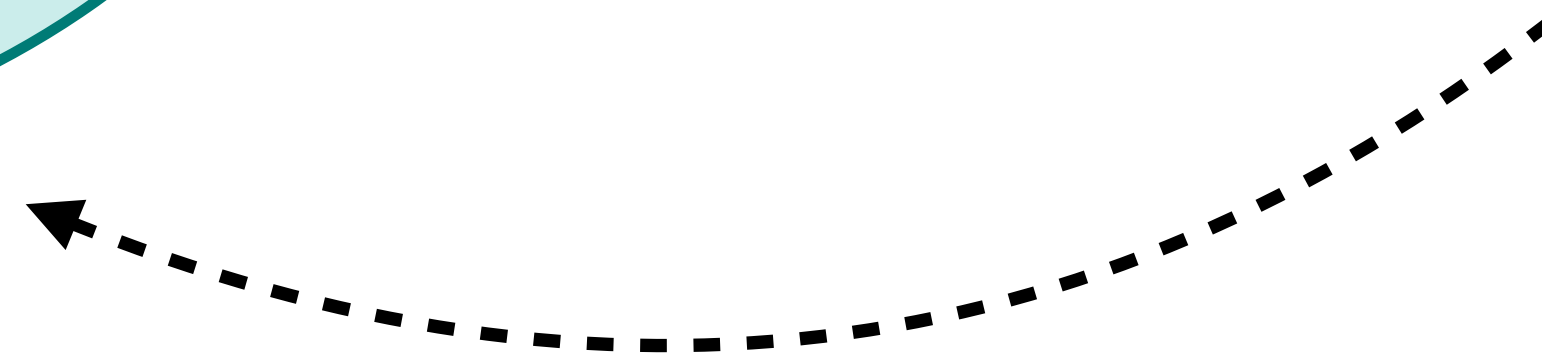
**Population**



**14x**

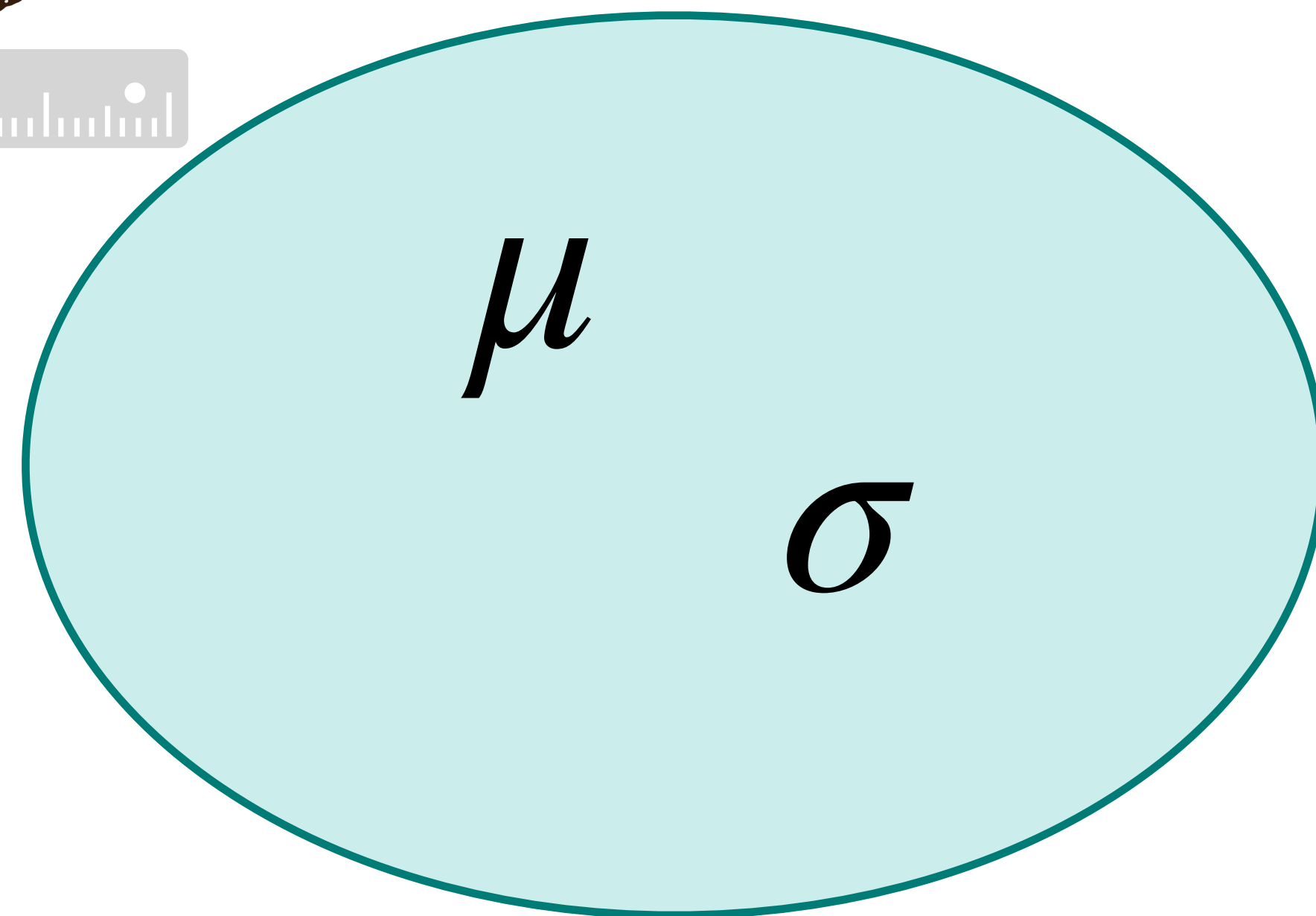


**Sample**

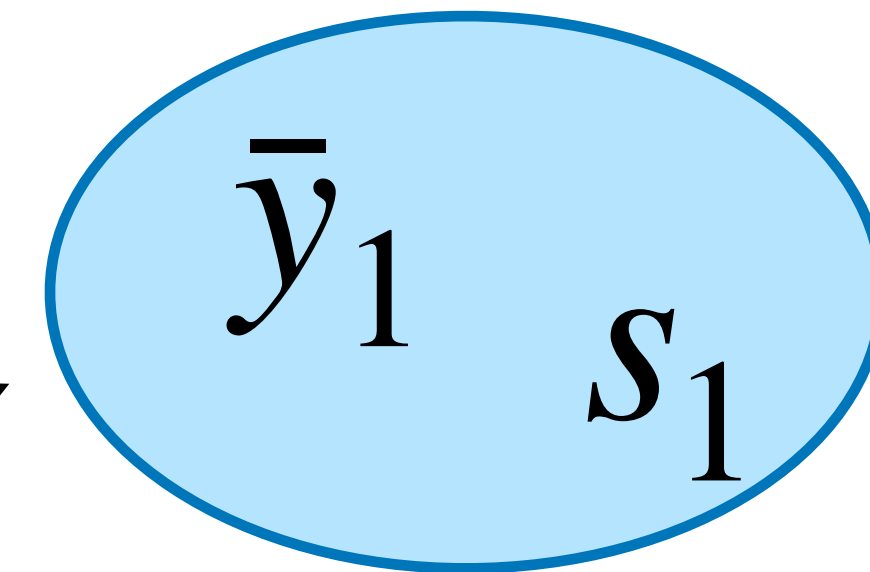
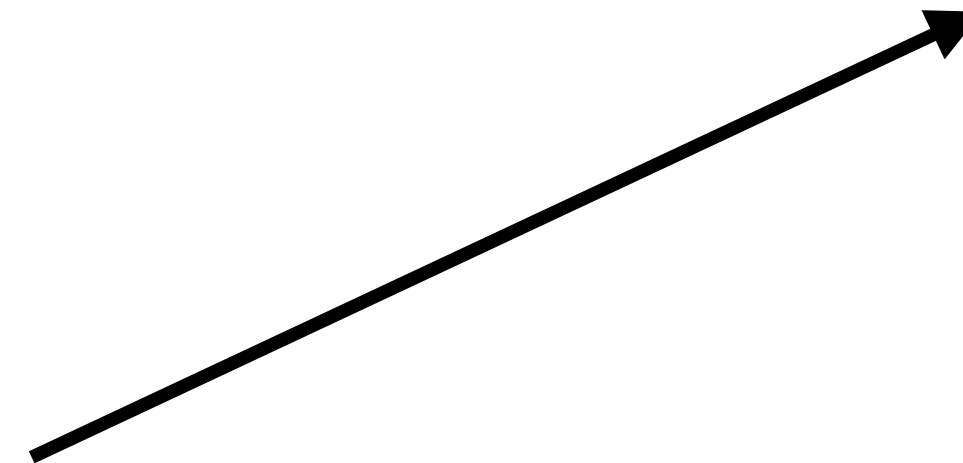




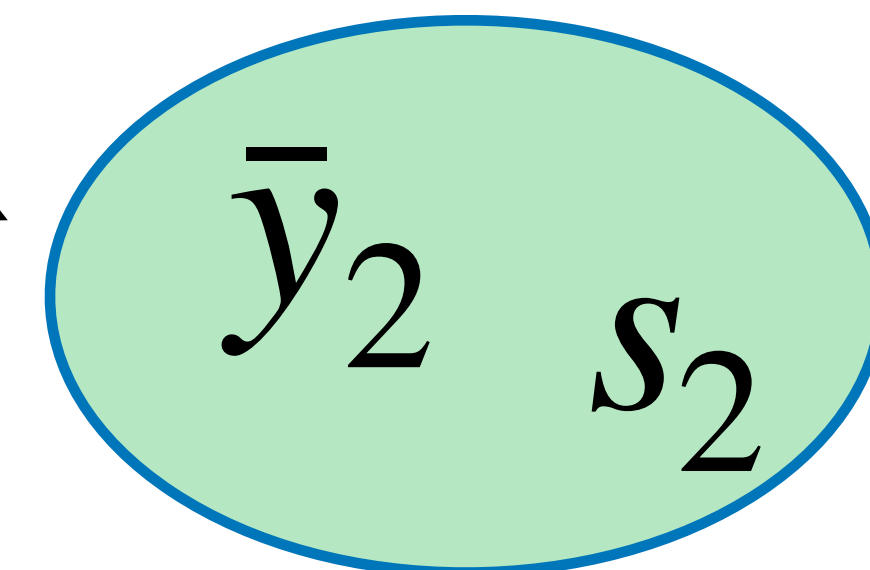
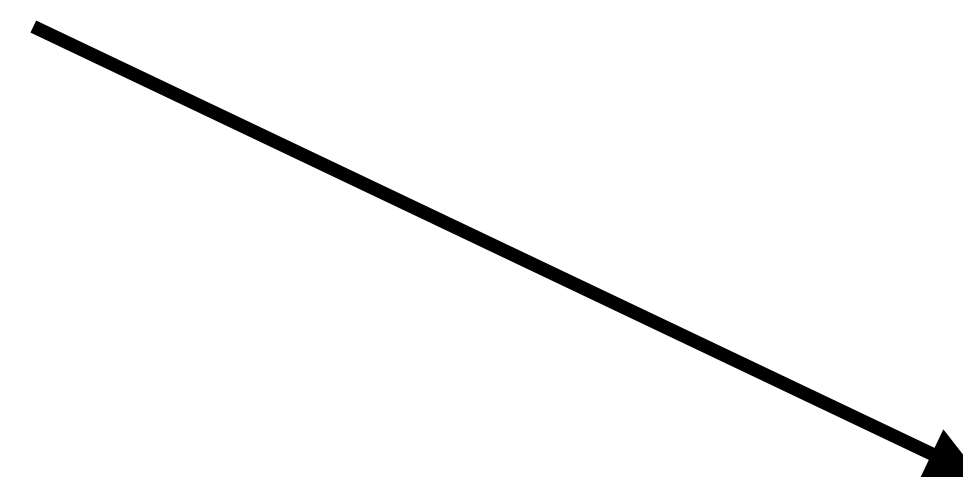
# Statistical estimation



**Population**



**Sample 1**

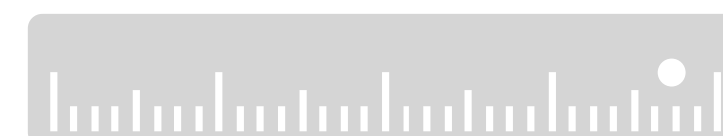


**Sample 2**

**14x**

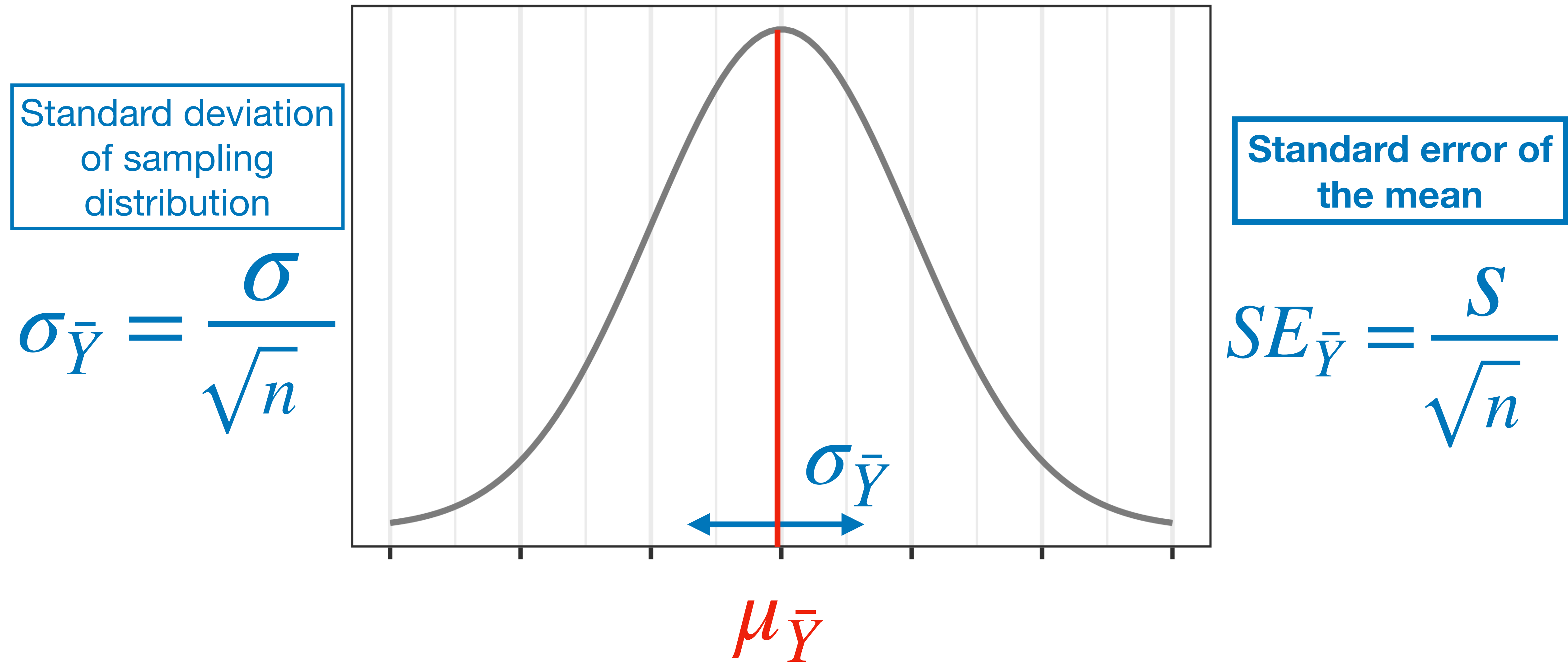


**14x**





**Standard error of the mean** is estimated from the sampling distribution



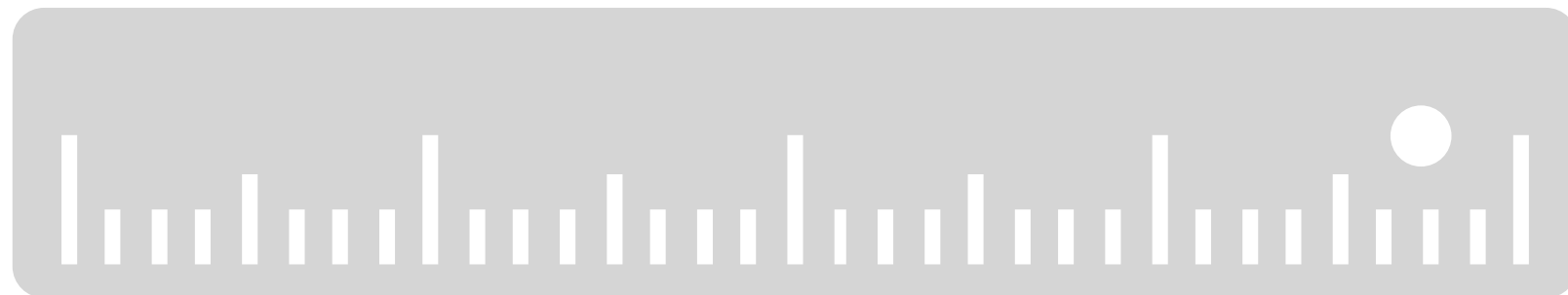
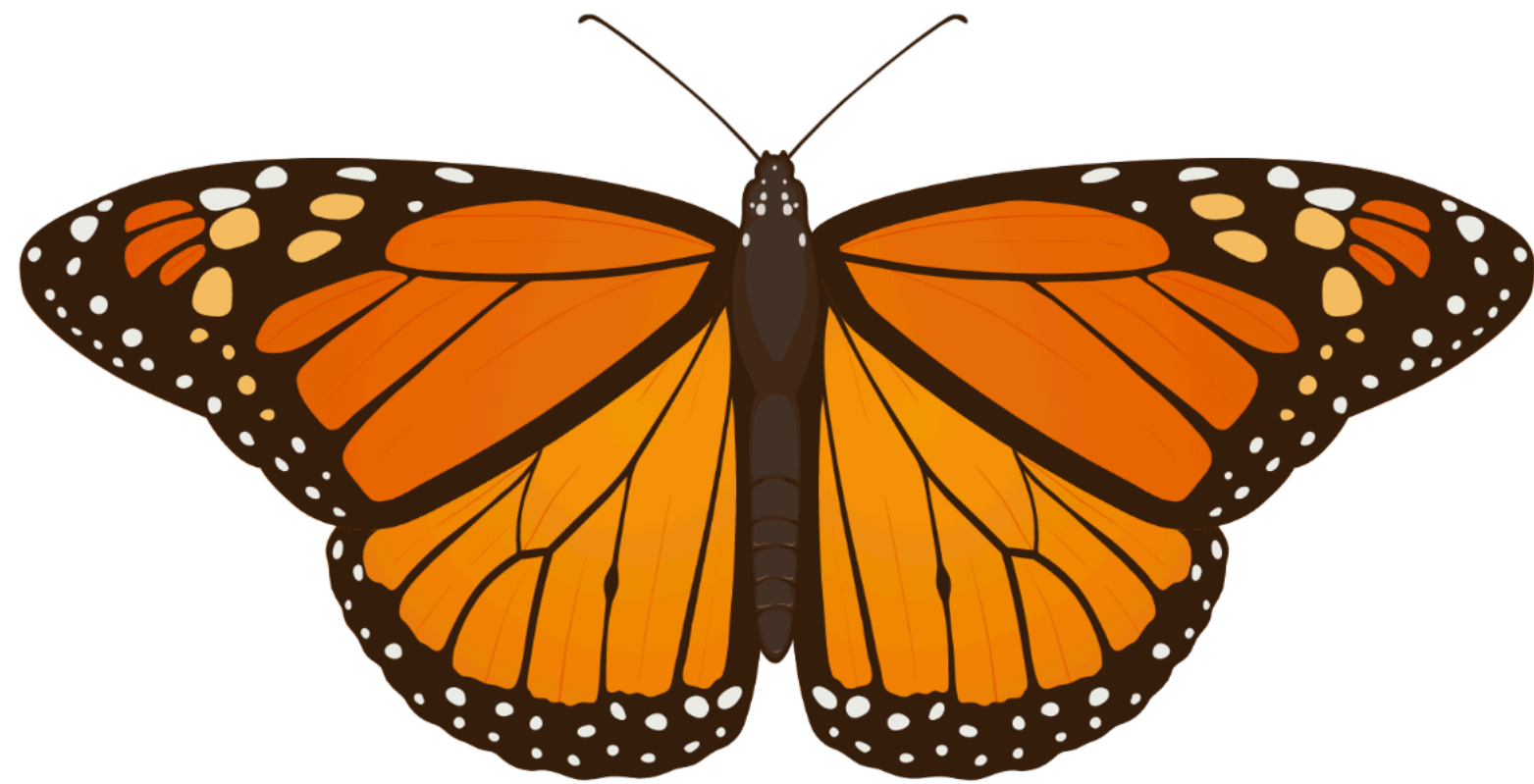
# Standard error of the mean is estimated from the sampling distribution

- **Standard error of the mean (SE)** is a measure of reliability or precision of the sample mean as an estimate of the population mean
- SE incorporates the two factors that influence reliability:
  - Variability of observations ( $s$ )
  - Sample size ( $n$ )

**Standard error of  
the mean**

$$SE_{\bar{Y}} = \frac{s}{\sqrt{n}}$$

# Standard error of the mean is estimated from the sampling distribution



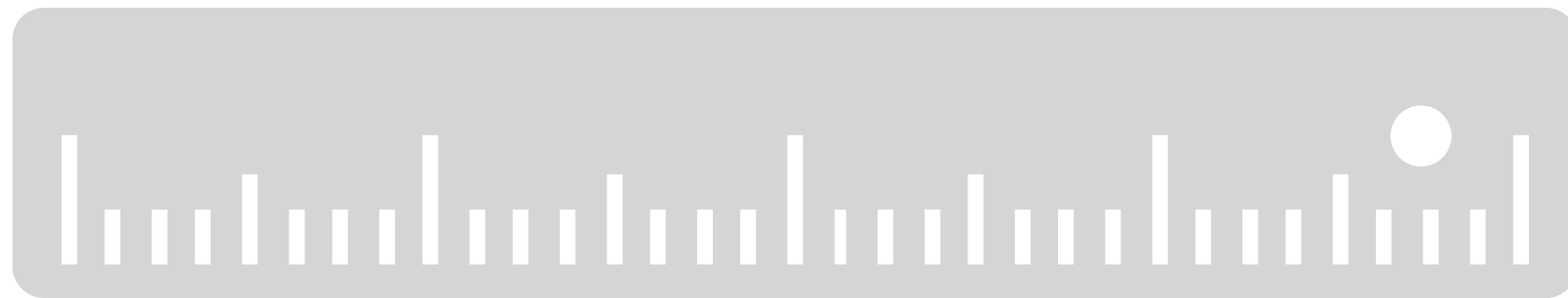
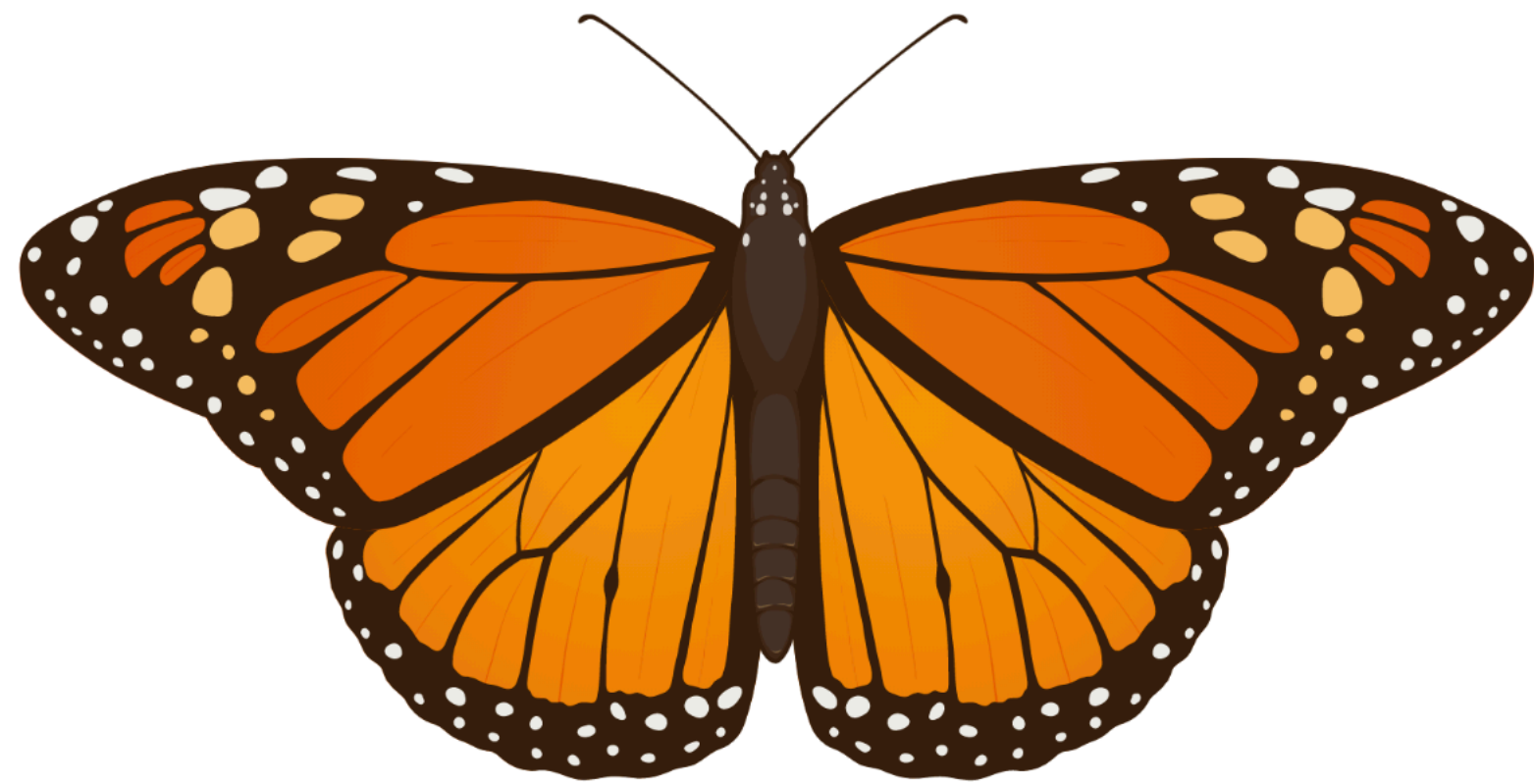
$n = 14$

$$\bar{y} = 32.81 \text{ cm}^2$$

$$s = 2.48 \text{ cm}^2$$

$$SE_{\bar{y}} = \frac{s}{\sqrt{n}}$$

# Standard error of the mean is estimated from the sampling distribution



$n = 14$



$$\bar{y} = 32.81 \text{ cm}^2$$

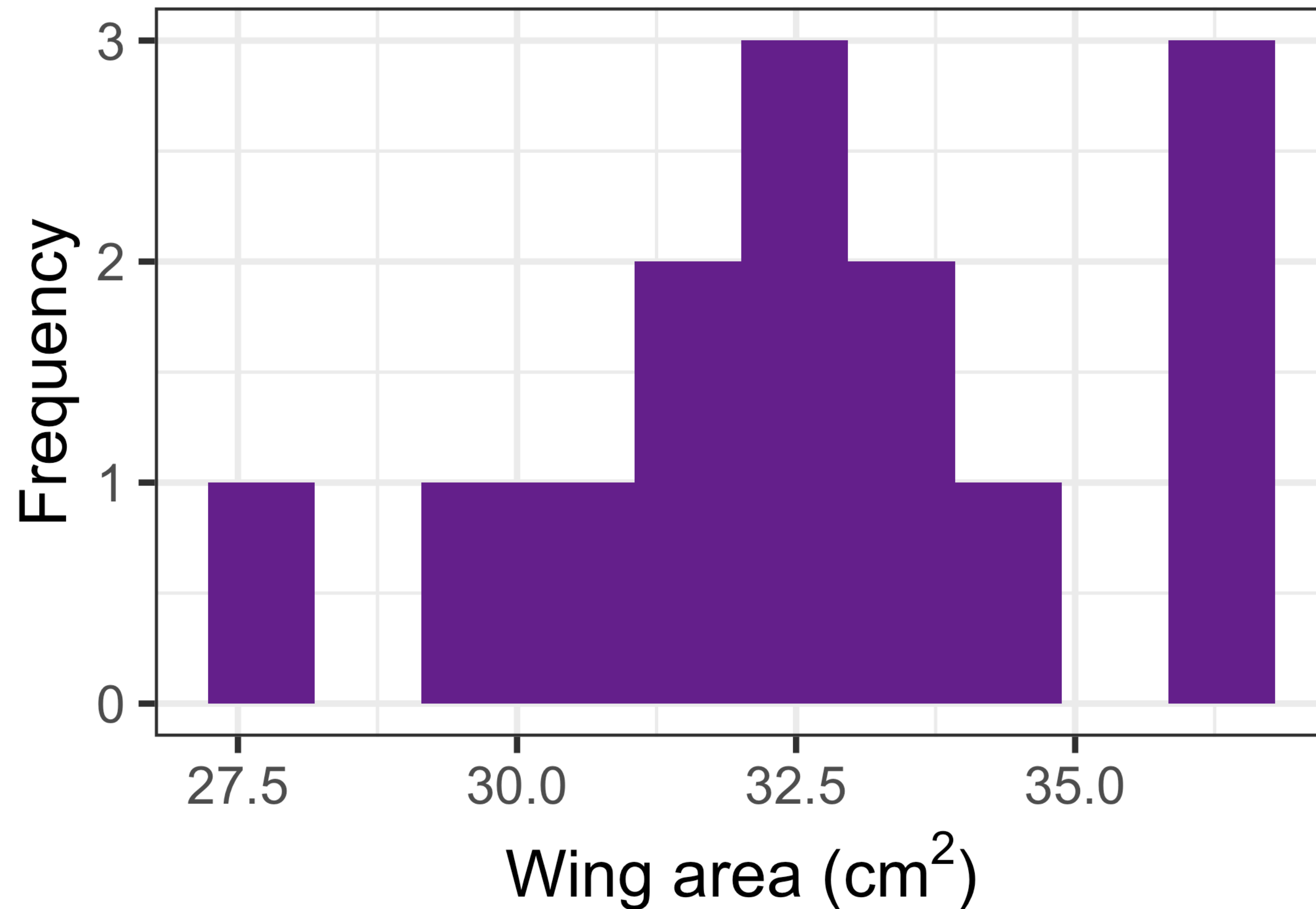
$$s = 2.48 \text{ cm}^2$$

$$SE_{\bar{y}} = \frac{2.48}{\sqrt{14}}$$

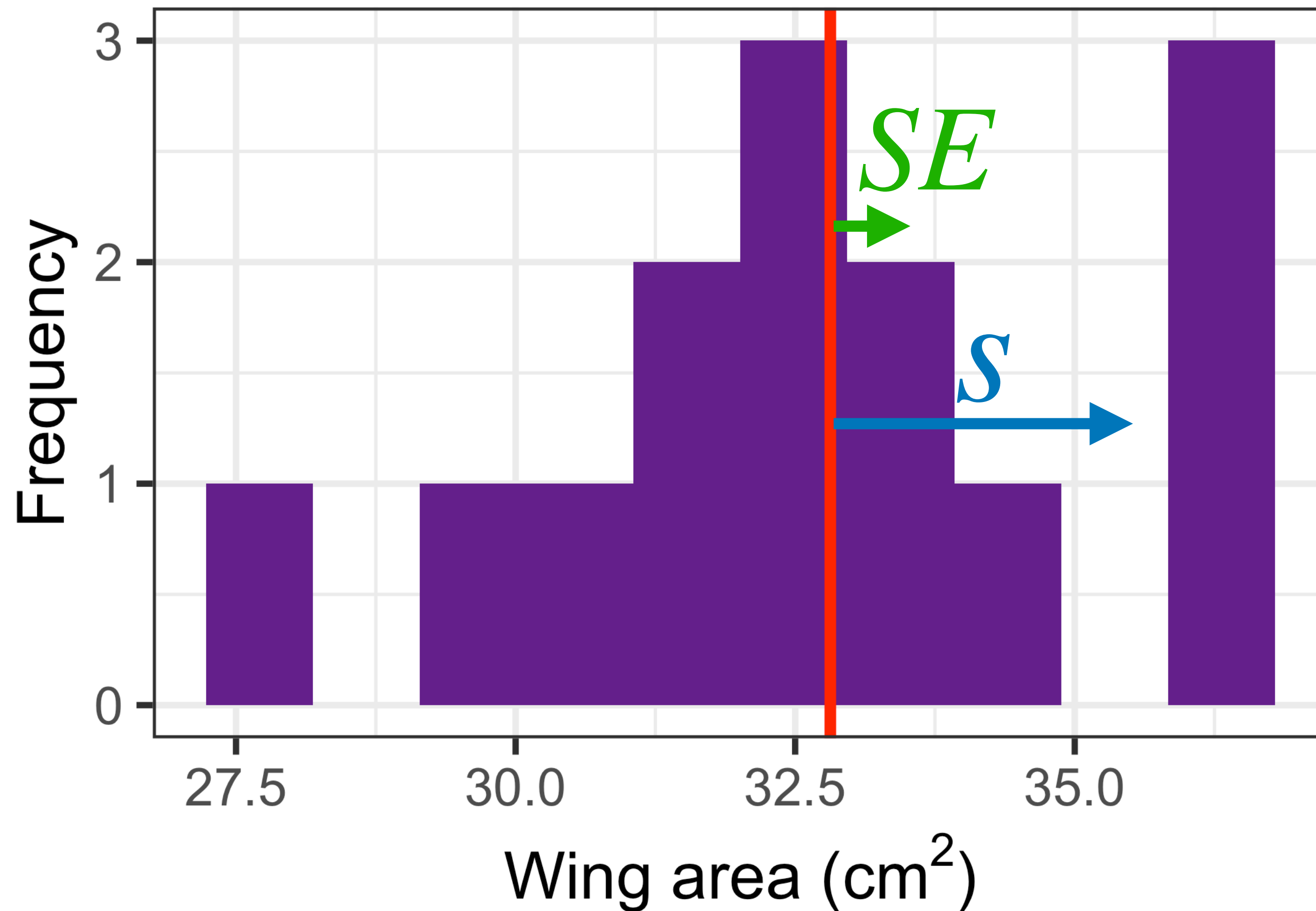
*What does the SE mean?*

$$= 0.66 \text{ cm}^2$$

# Standard error of the mean is estimated from the sampling distribution

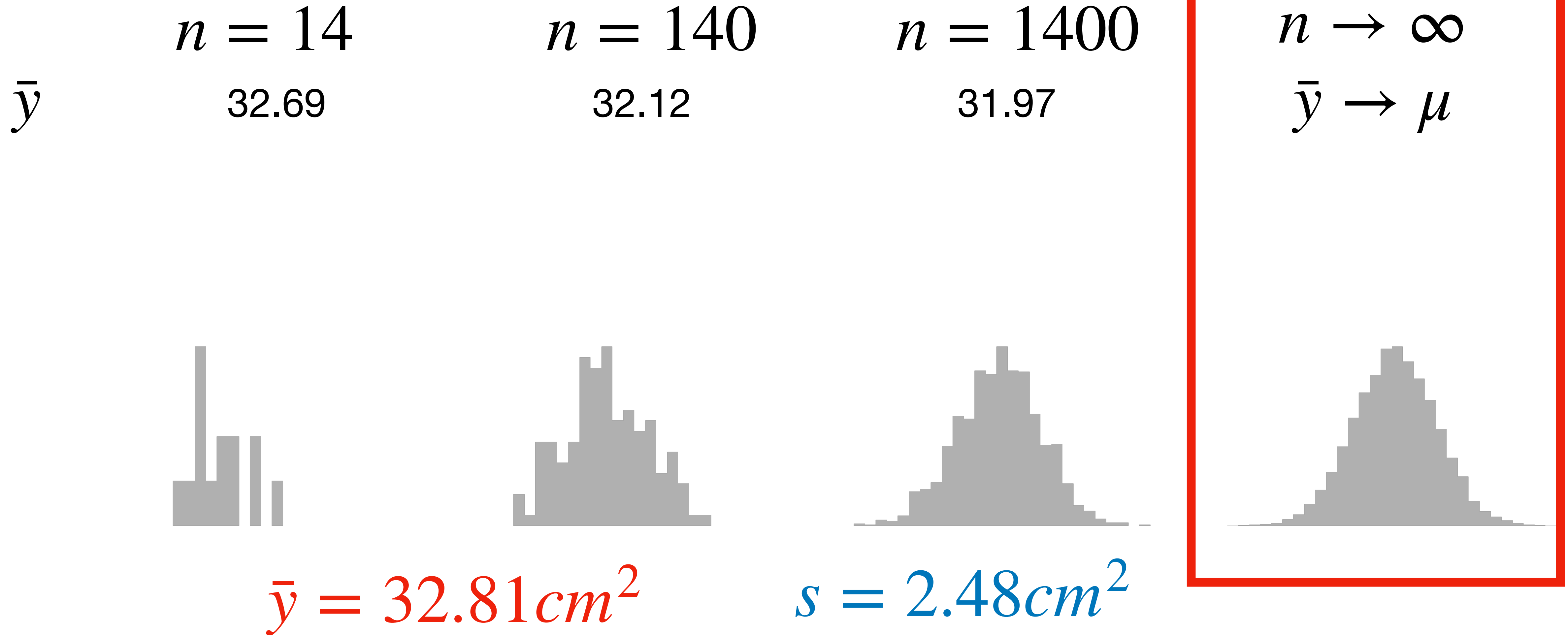


# Standard error of the mean is estimated from the sampling distribution

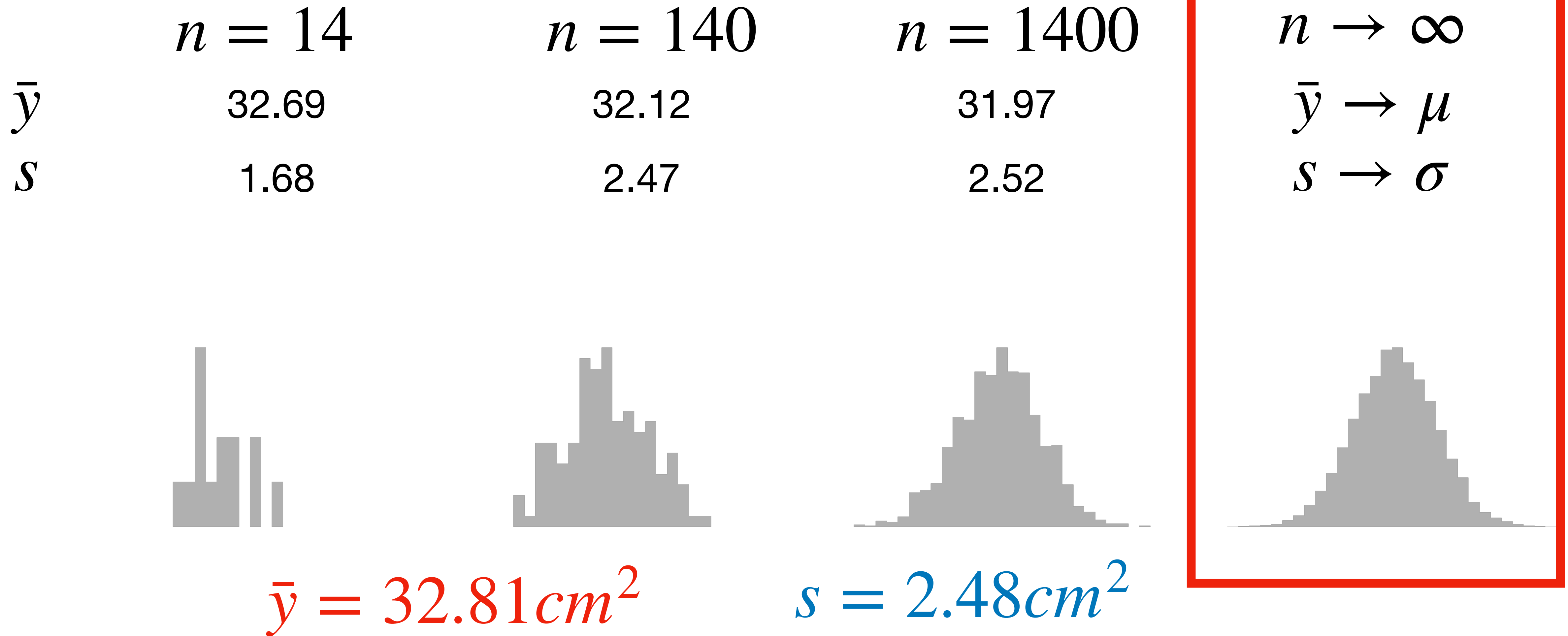




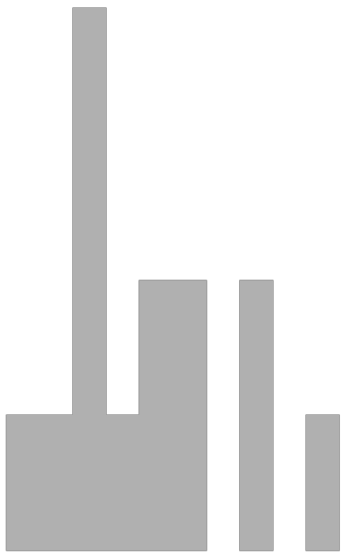


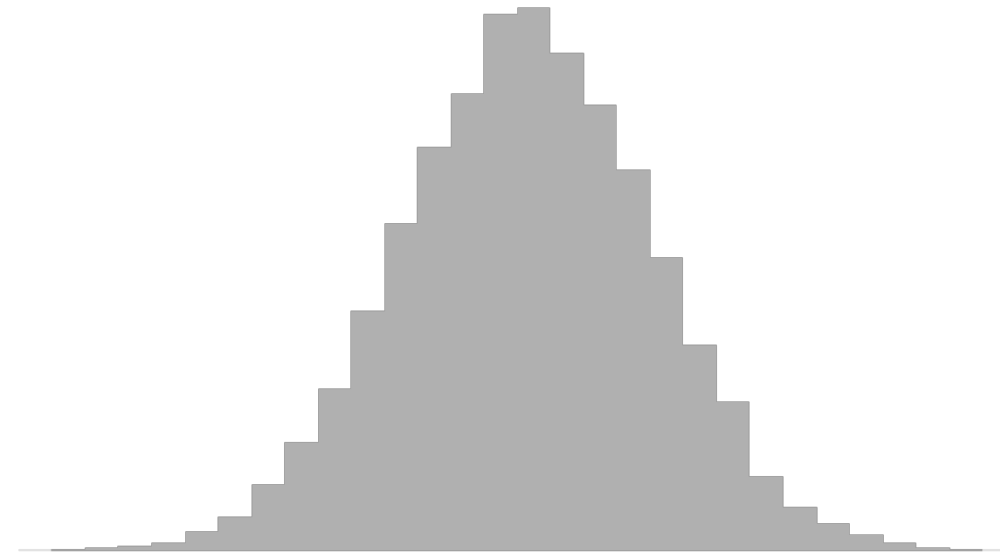
# Standard error (SE) versus standard deviation (SD)



# Standard error (SE) versus standard deviation (SD)

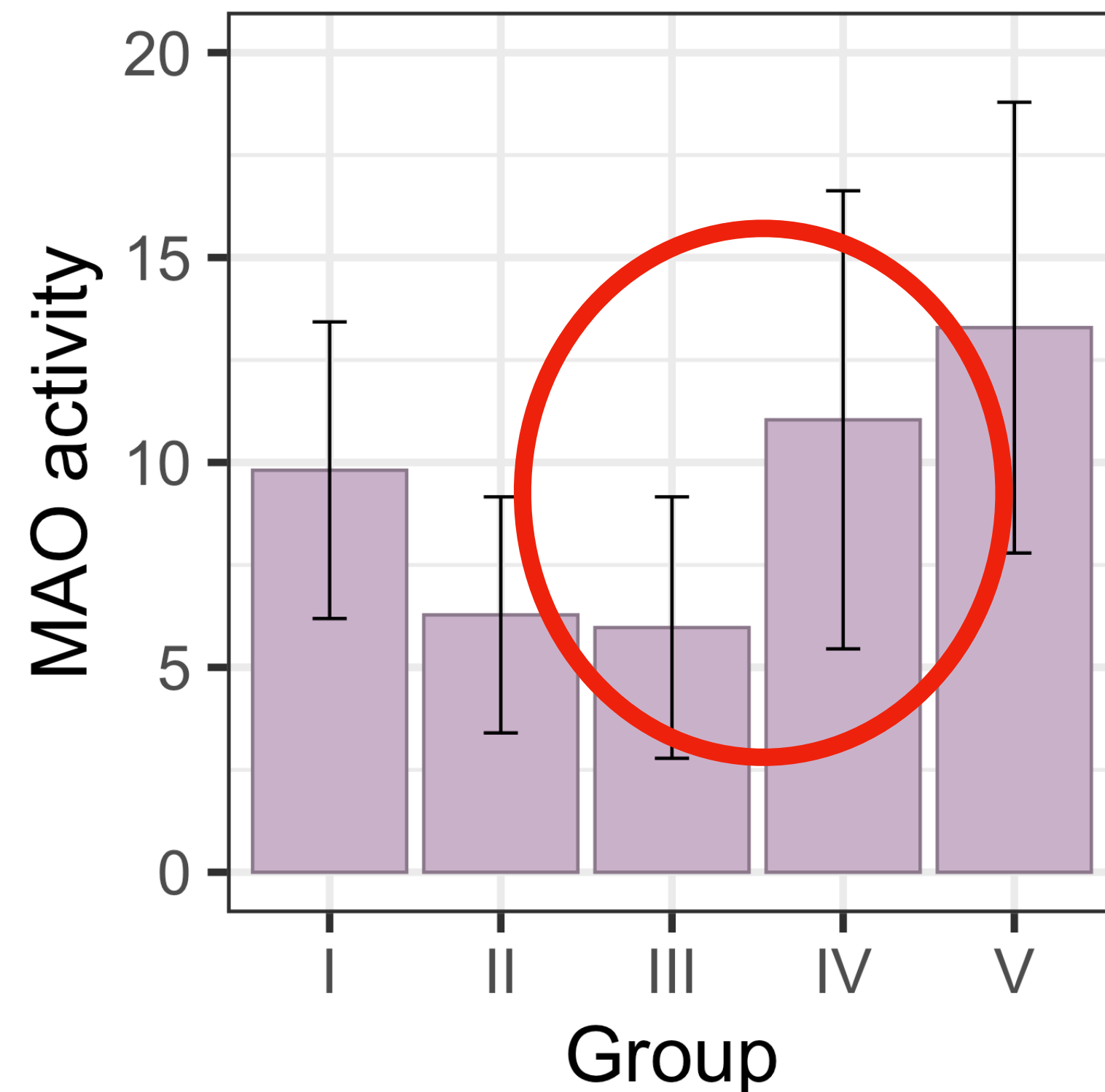


# Standard error (SE) versus standard deviation (SD)

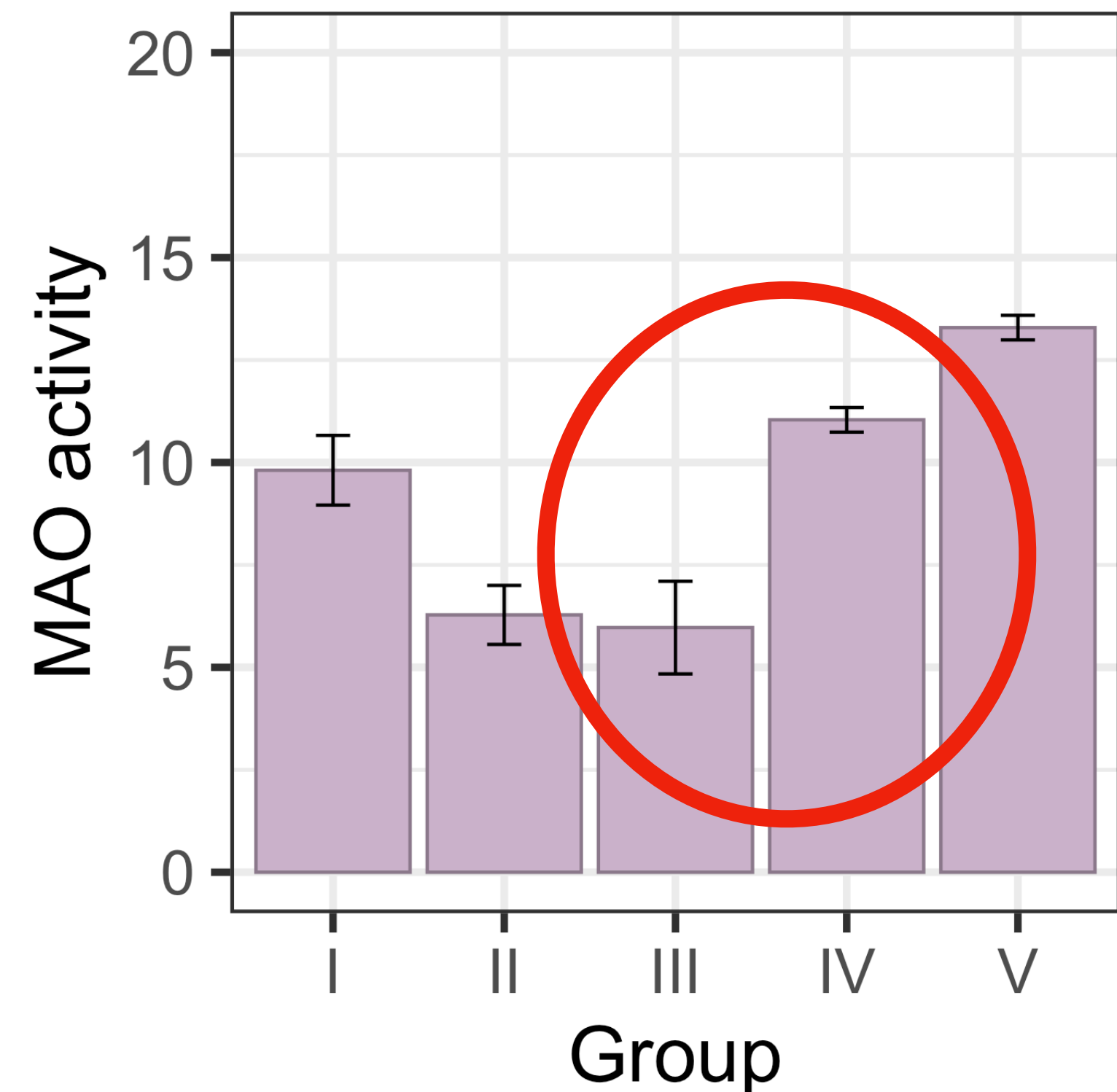
	$n = 14$	$n = 140$	$n = 1400$	$n \rightarrow \infty$
$\bar{y}$	32.69	32.12	31.97	$\bar{y} \rightarrow \mu$
$s$	1.68	2.47	2.52	$s \rightarrow \sigma$
$SE$	0.451	0.209	0.067	$SE \rightarrow 0$
				
	$\bar{y} = 32.81 cm^2$		$s = 2.48 cm^2$	

# Standard error (SE) versus standard deviation (SD)

**SD = dispersion of data**



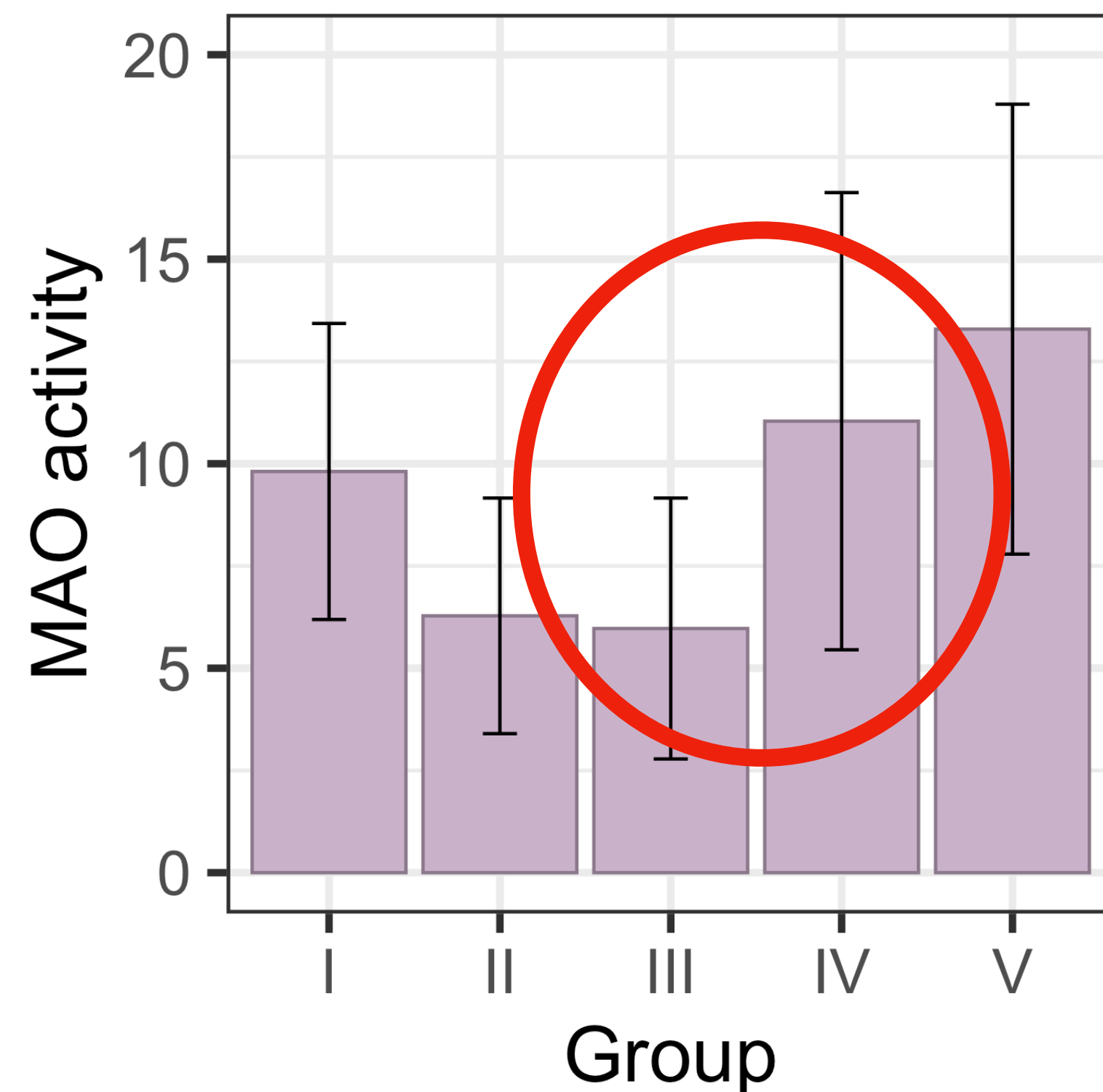
**SE = unreliability in the estimate  
of the population mean**



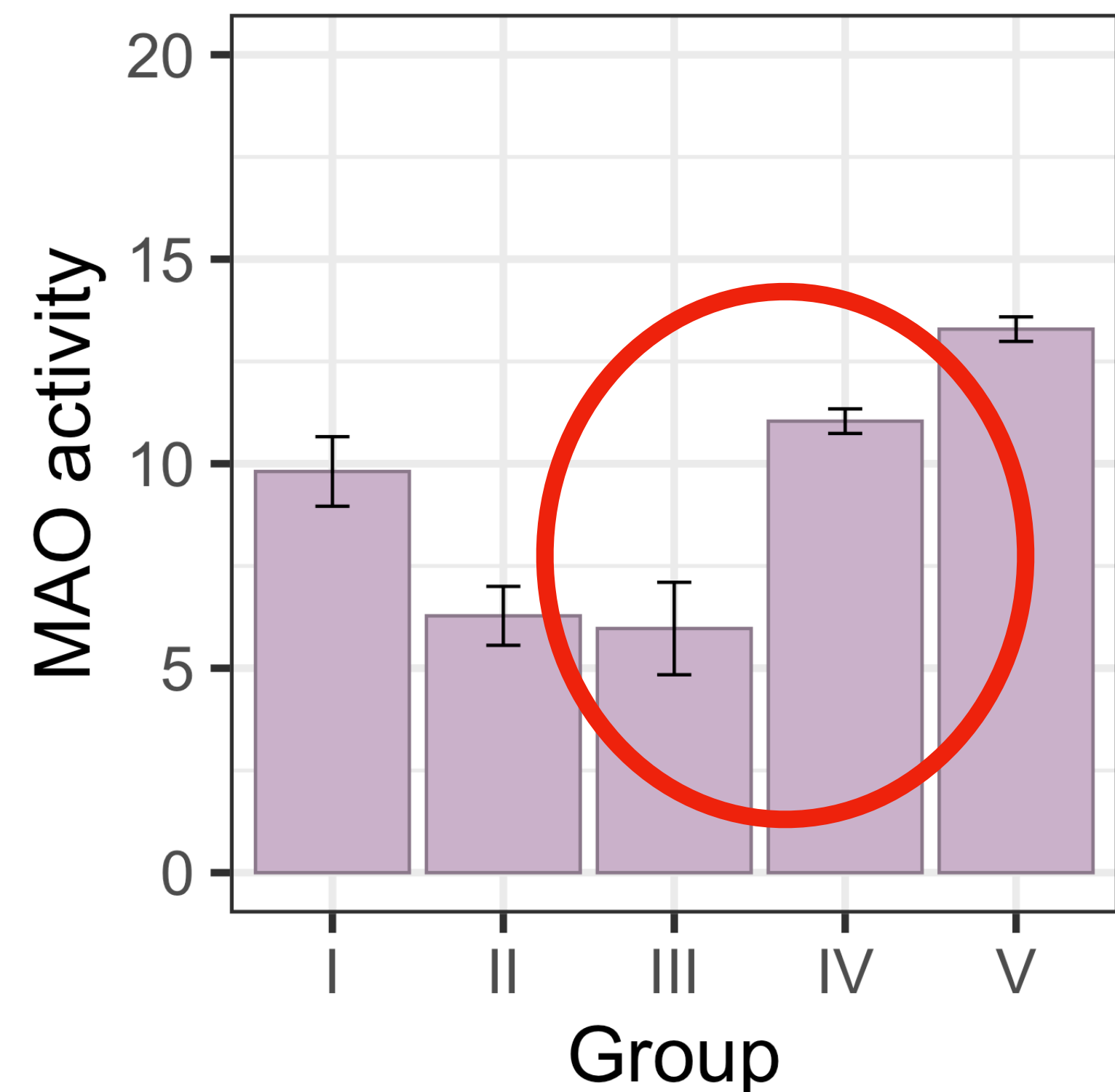
# When would we choose to plot SE v. SD?

*Do you want to compare means or summarize data variability?*

**SD = dispersion of data**



**SE = unreliability in the estimate of the population mean**



# Example

**A pharmacologist measured the concentration of dopamine in the brains of eight rats. The mean concentration was 1,269 ng/gm and the standard deviation was 145 ng/gm. What was the standard error of the mean?**

$$SE = \frac{s}{\sqrt{n}}$$

$$SE = \frac{145}{\sqrt{8}} = 51.2$$



# Practice

**This quantity is a measure of the accuracy of the sample mean as an estimate of the population mean**

Standard Error (SE)

Standard Deviation (SD)

**This quantity tends to stay the same as the sample size goes up**

Standard Error (SE)

Standard Deviation (SD)

**This quantity tends to go down as the sample size goes up**

Standard Error (SE)

Standard Deviation (SD)

# Practice

**This quantity is a measure of the accuracy of the sample mean as an estimate of the population mean**

Standard Error (SE)

Standard Deviation (SD)

**This quantity tends to stay the same as the sample size goes up**

Standard Error (SE)

Standard Deviation (SD)

**This quantity tends to go down as the sample size goes up**

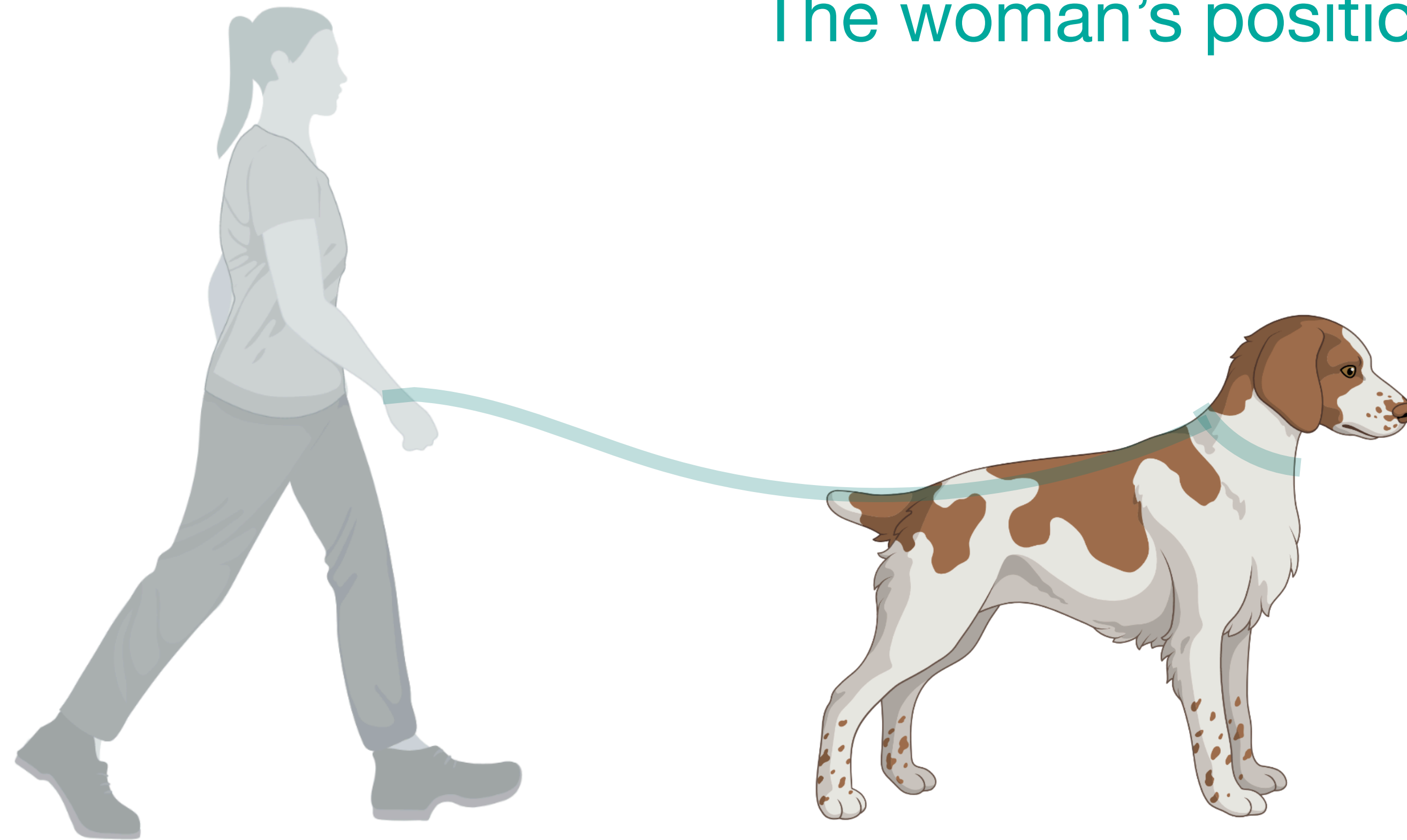
Standard Error (SE)

Standard Deviation (SD)

# The confidence interval



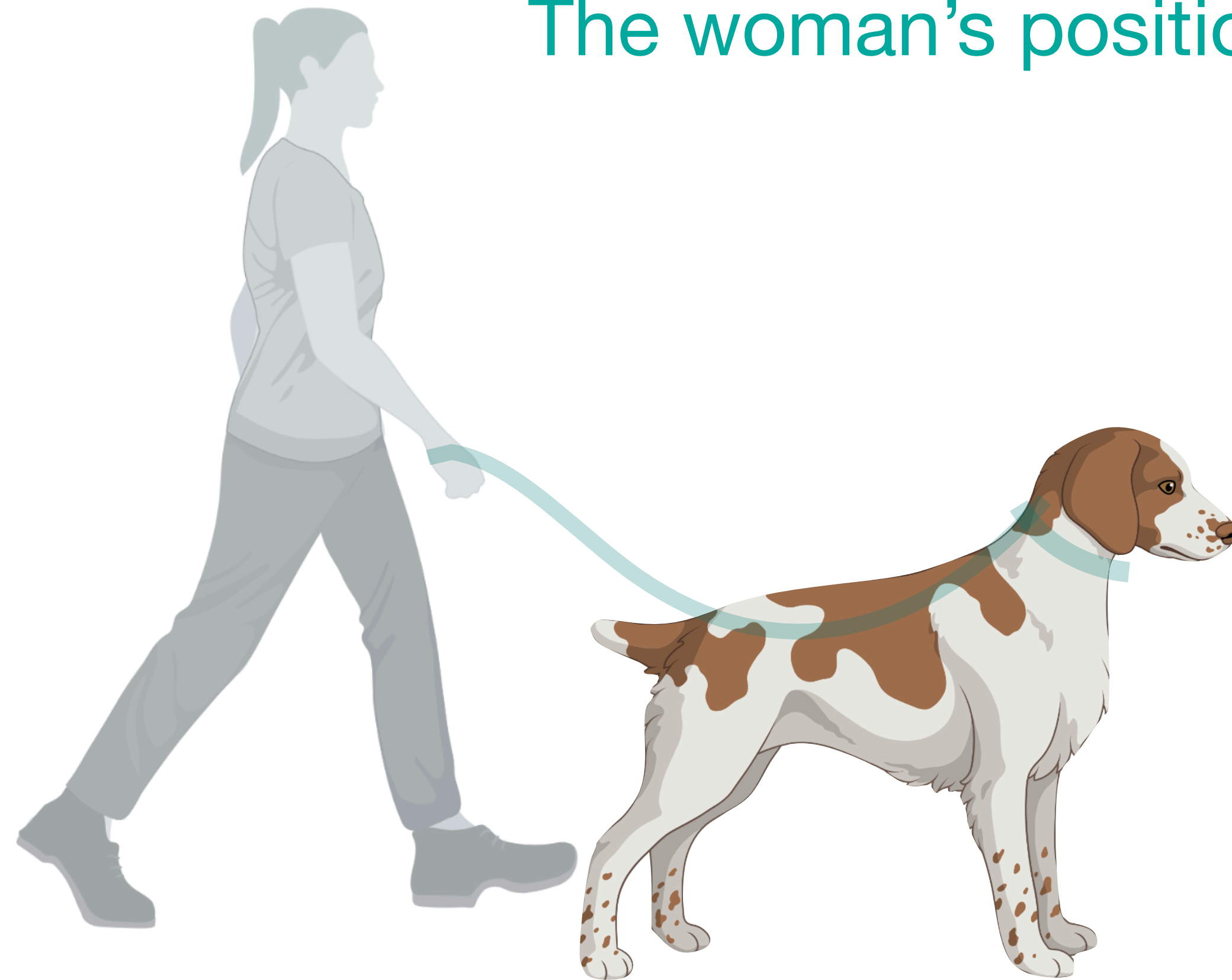
The woman's position



# The confidence interval



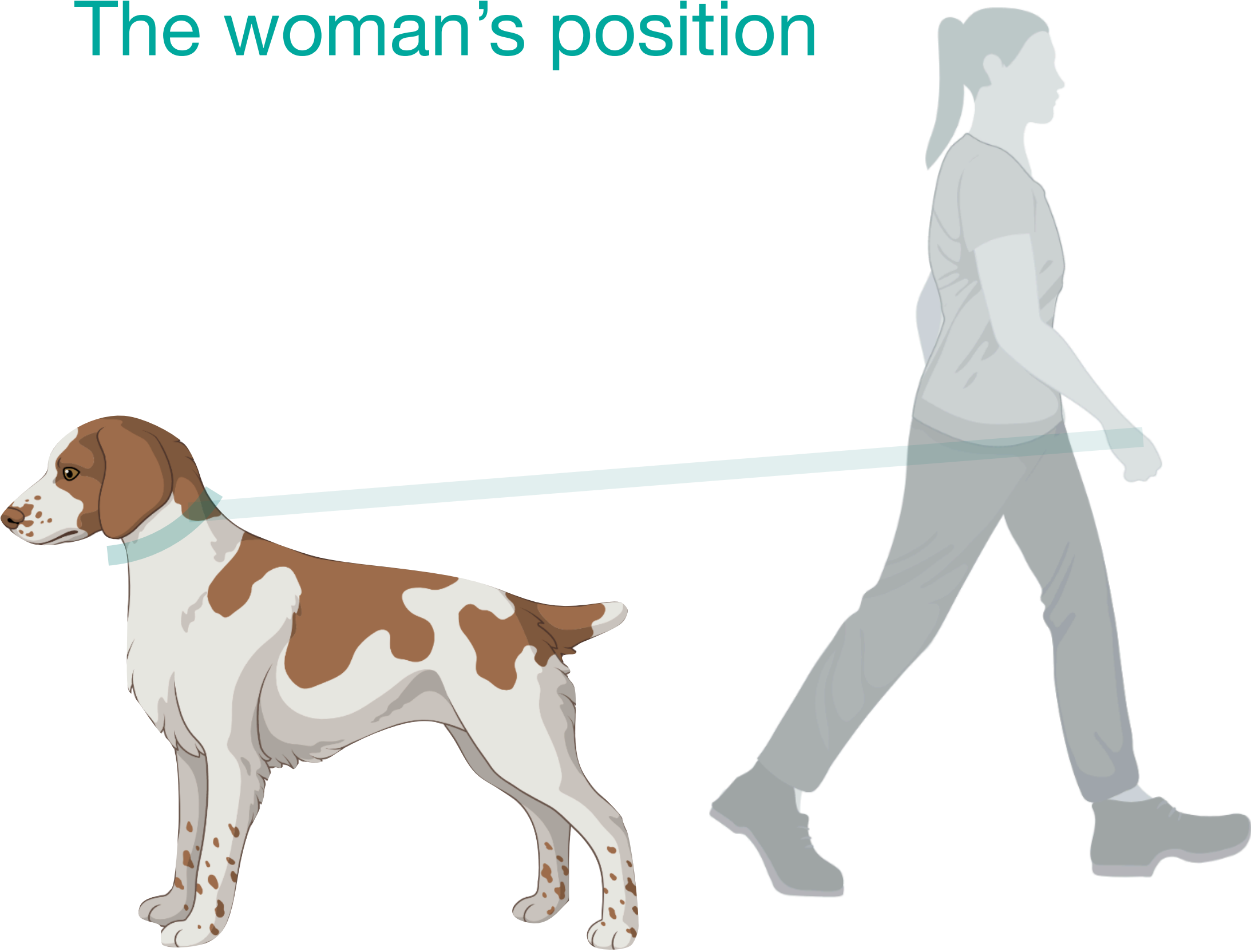
The woman's position



# The confidence interval



The woman's position

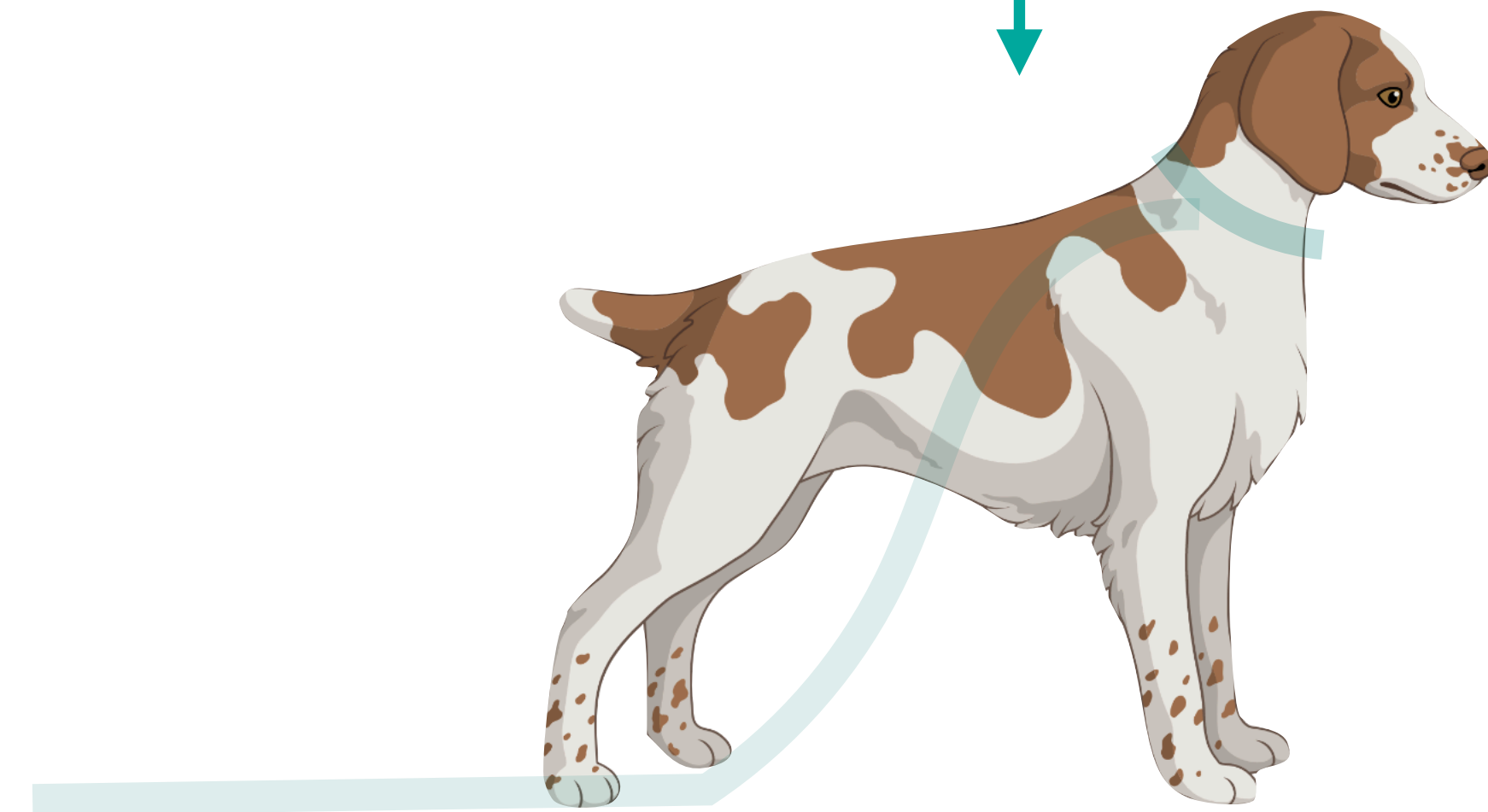


# The confidence interval



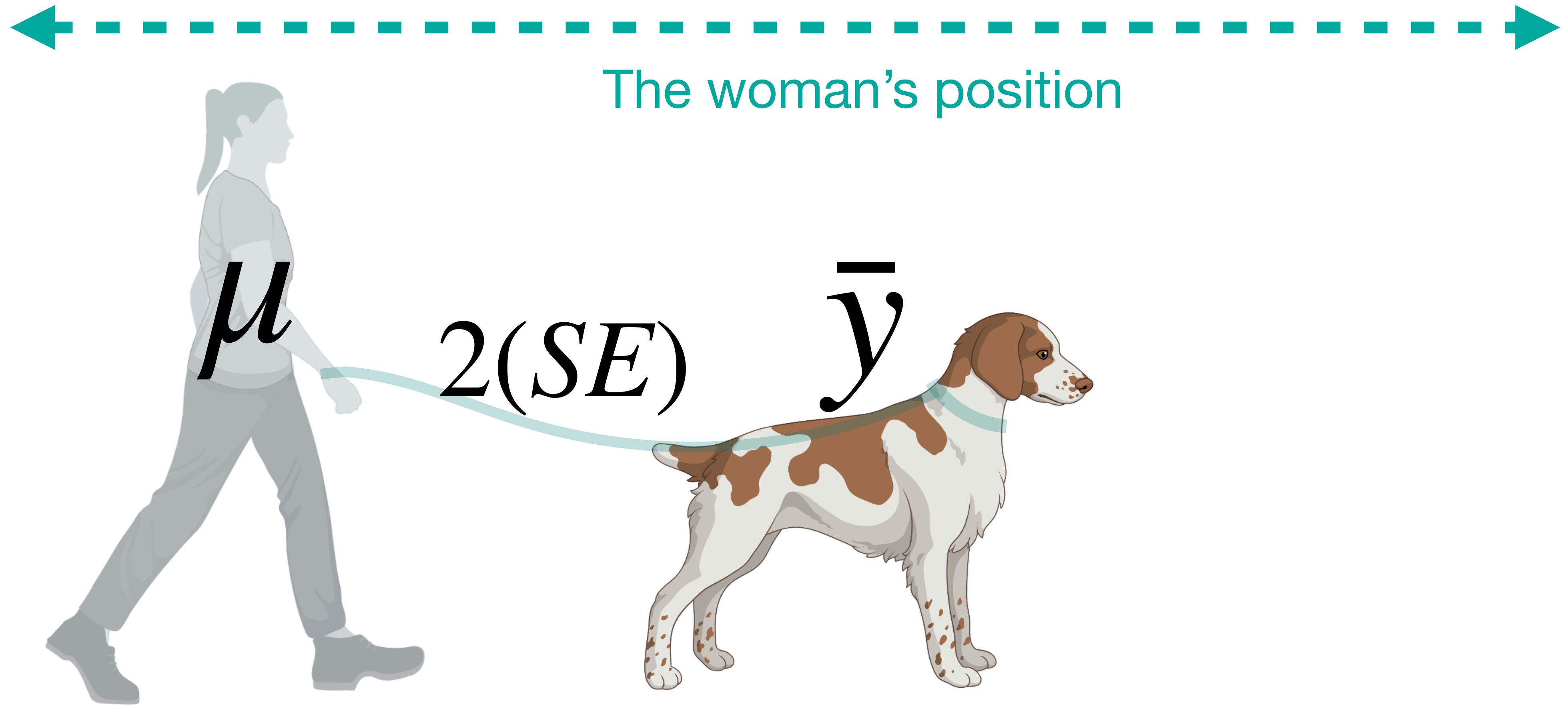
The woman's position

*Not likely, but possible*



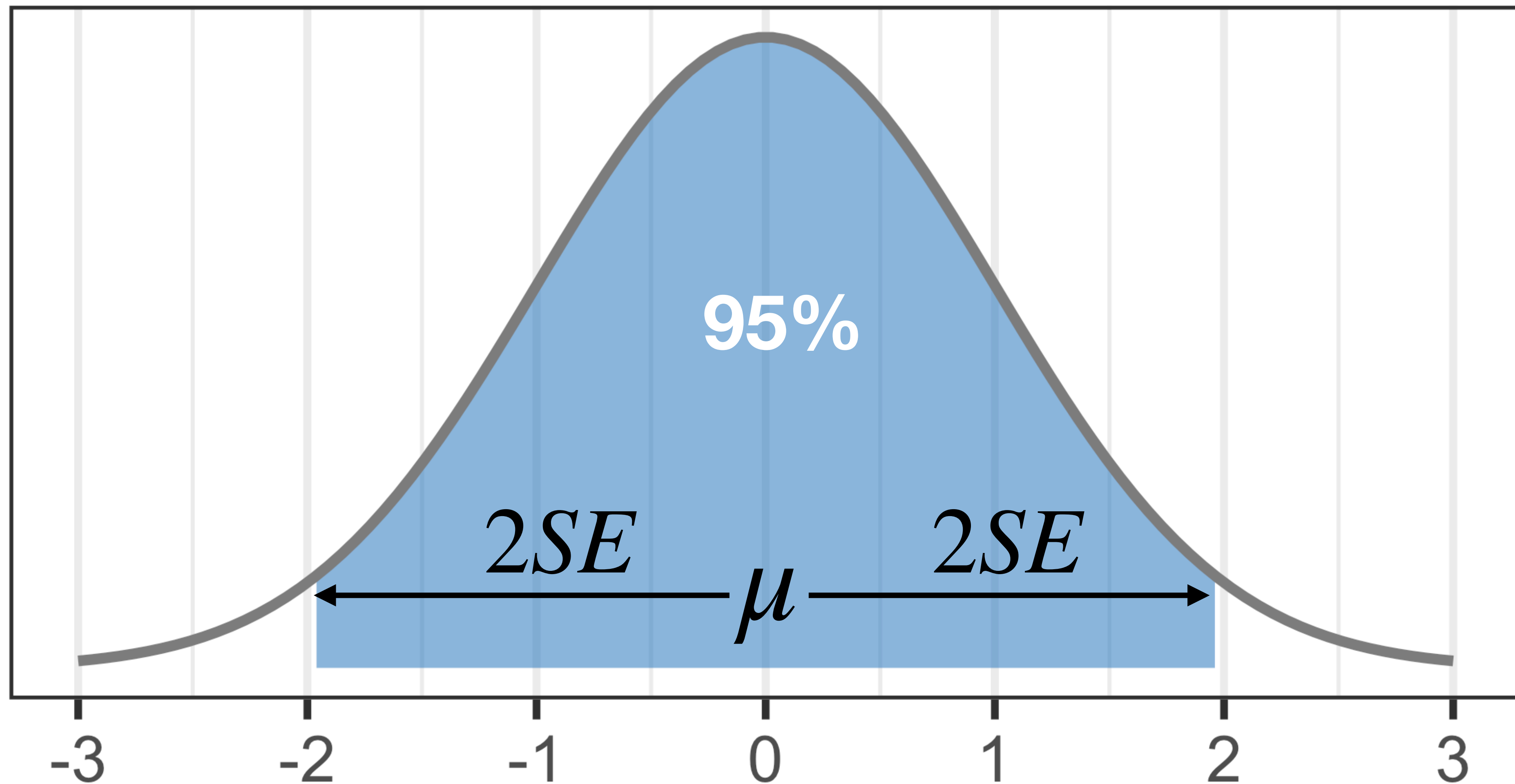


# The confidence interval



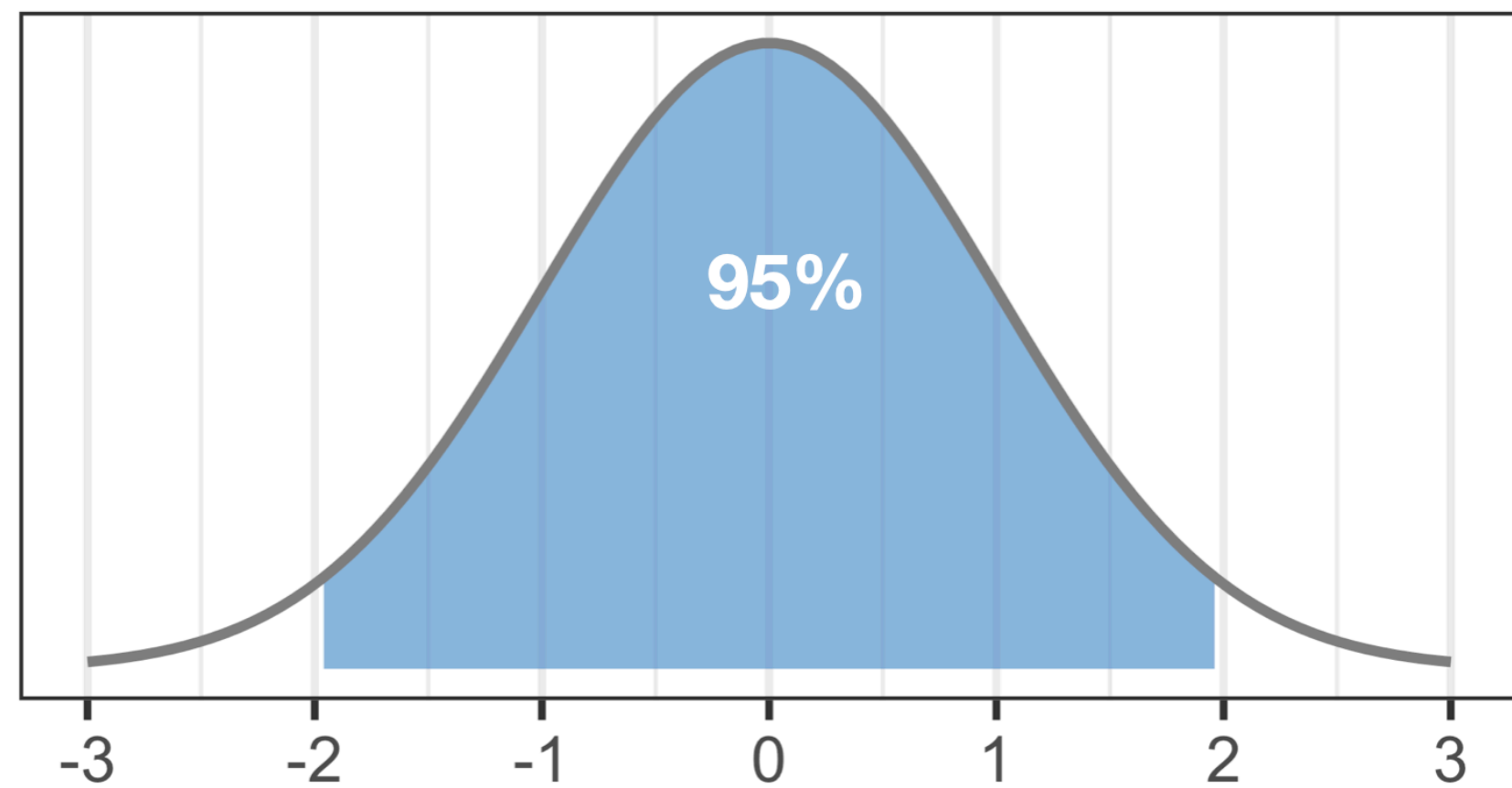
**95% confidence interval of woman's position = position of the dog  $\pm 2 \times SE$**

# The confidence interval



**Sampling distribution of Y - random sample from normal distribution**

# The confidence interval



Will contain  $\mu$   
for 95% of all samples



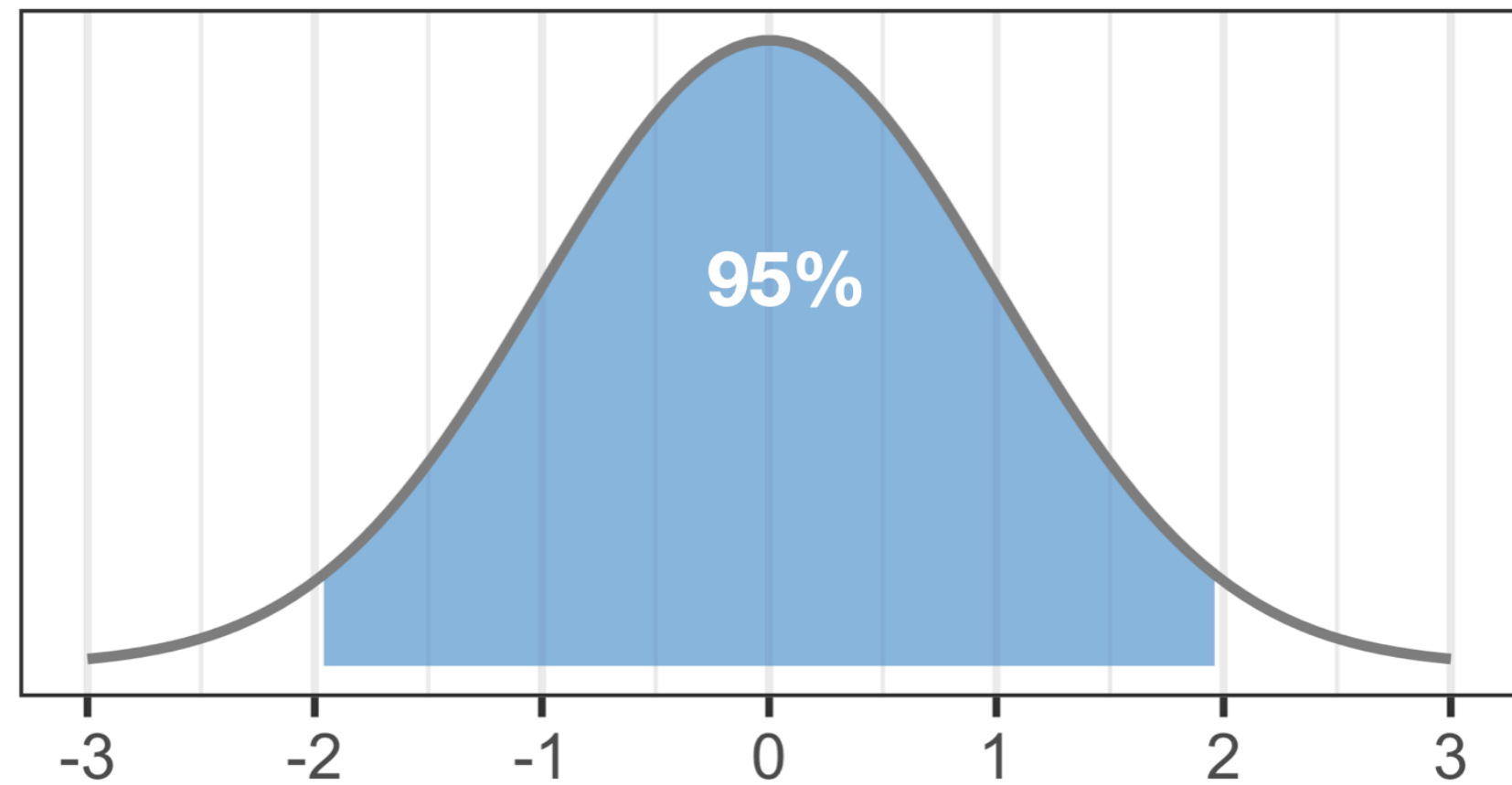
$$Pr[-1.96 < Z < 1.96] = 0.95$$

$$\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}$$

But... we don't  
know  $\sigma$

$$\bar{Y} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$

# Student's $t$ distribution for confidence intervals



“Critical  
value”

Will contain  $\mu$   
for 95% of all samples

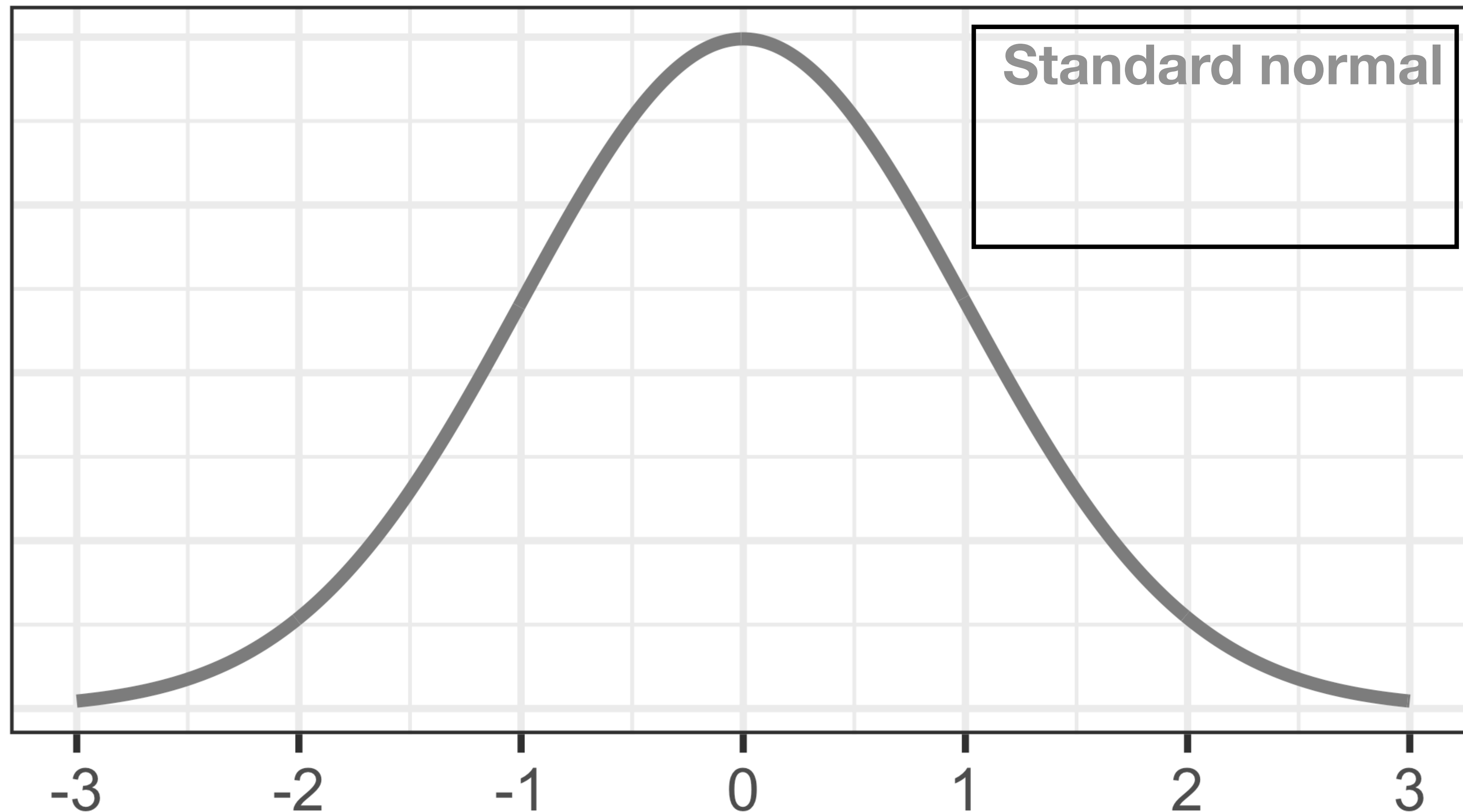
$$Pr[-1.96 < Z < 1.96] = 0.95$$

$$\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}$$

But... we don't  
know  $\sigma$

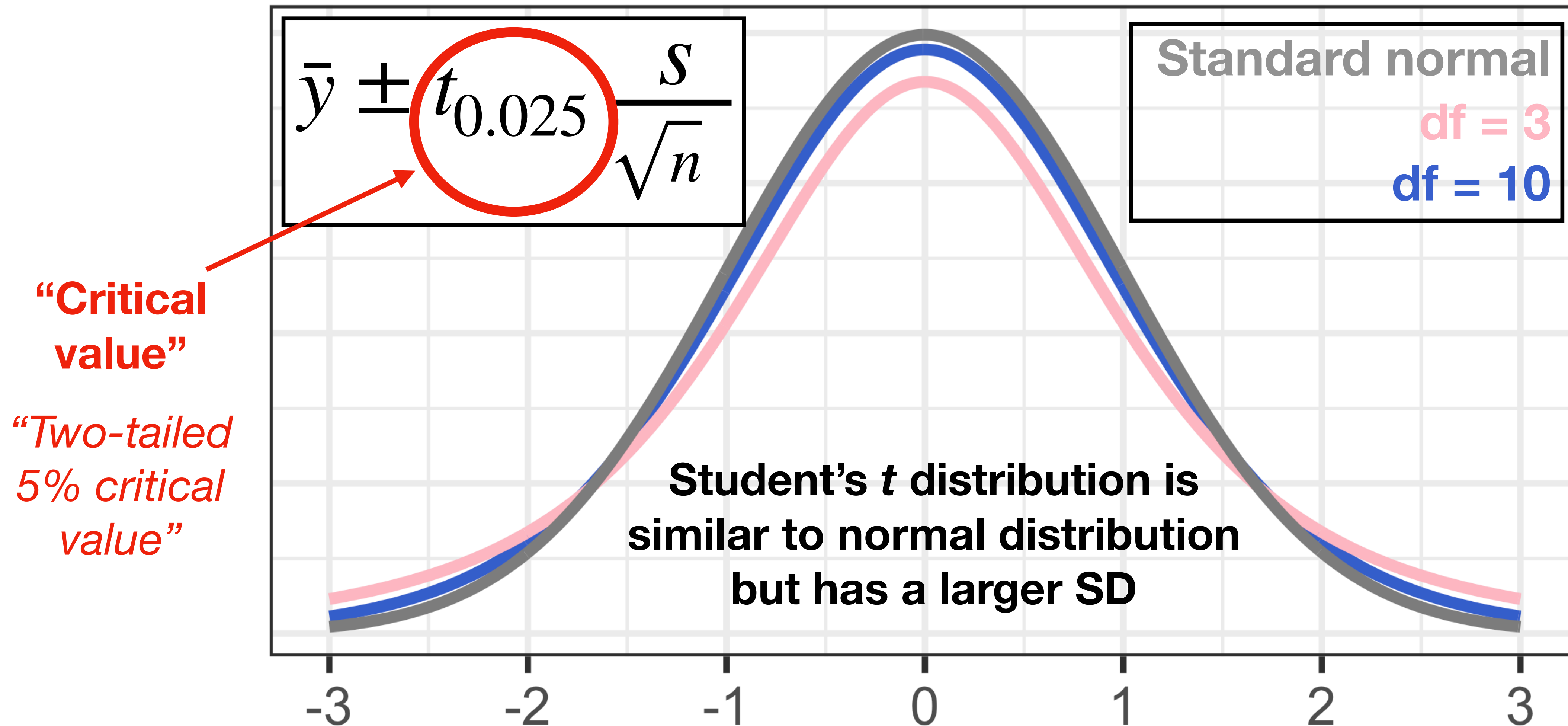
$$\bar{y} \pm t_{0.025} \frac{s}{\sqrt{n}}$$

# Student's $t$ distribution for confidence intervals



# Student's $t$ distribution for confidence intervals

Shape of distribution depends on **degrees of freedom**  $\longrightarrow$  ( $df = n - 1$ )

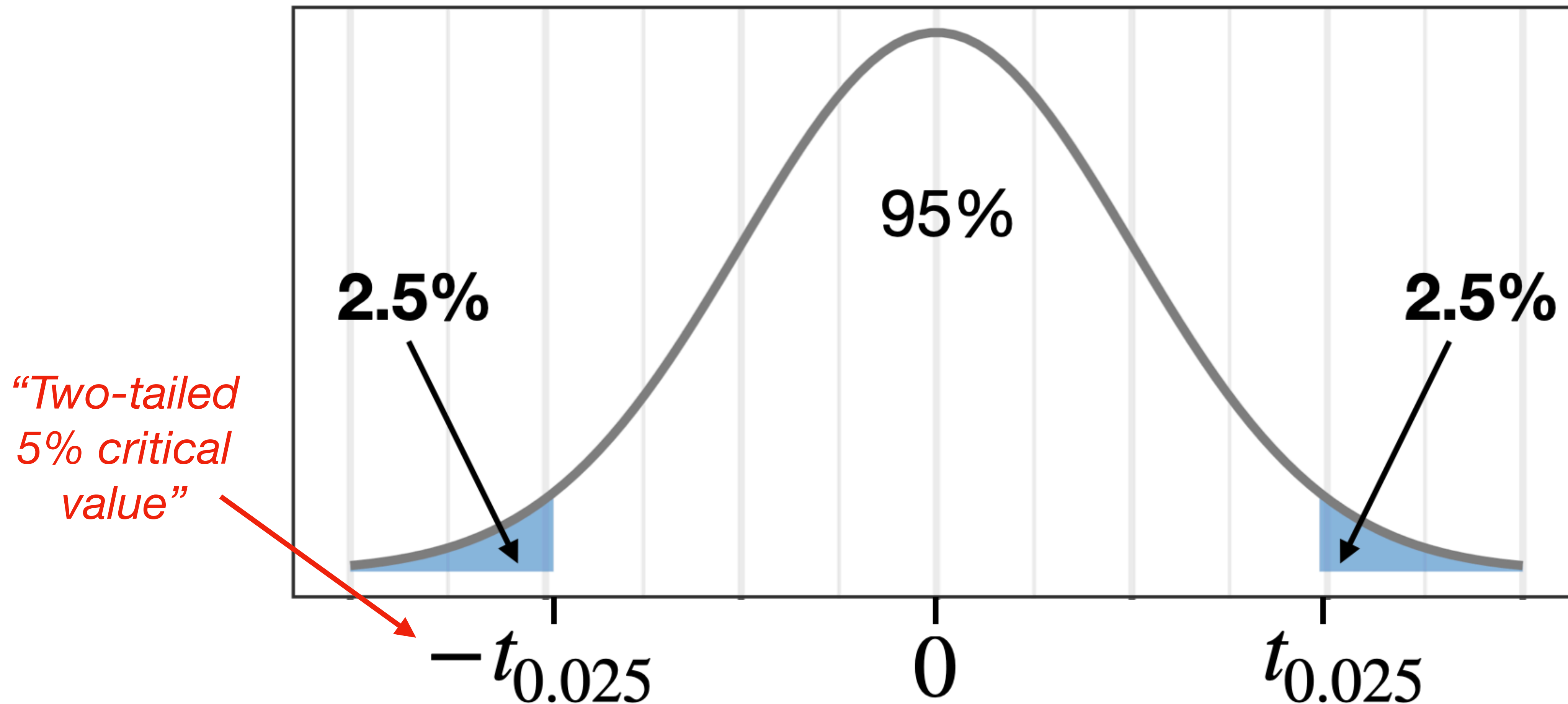




# Critical value and Student's $t$ distribution

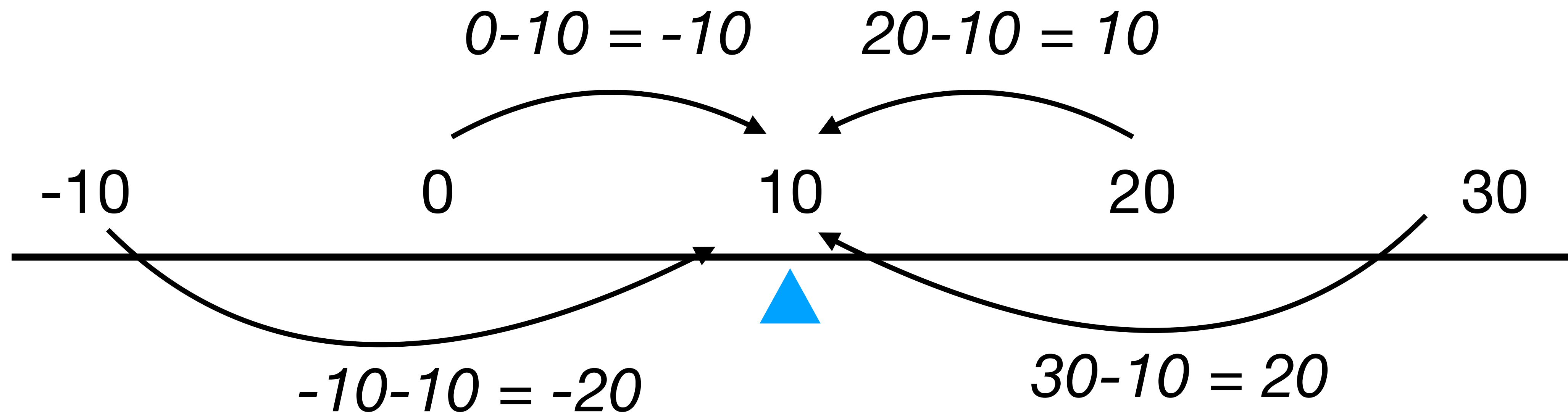
***“Two-tailed 5% critical value” = Combined area above  $t$  and below  $-t$  = 5%***

***“Two-tailed 5% critical value” = Area between two tails = 95%***



	P						
one-tail	0.1	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	0.2	0.1	0.05	0.02	0.01	0.002	0.001
DF							
1	3.078	6.314	12.706	31.821	63.656	318.289	636.578
2	1.886	2.92	4.303	6.965	9.925	22.328	31.6
3	1.638	2.353	3.182	4.541	5.841	10.214	12.924
4	1.533	2.132	2.776	3.747	4.604	7.173	8.61
5	1.476	2.015	2.571	3.365	4.032	5.894	6.869
6	1.44	1.943	2.447	3.143	3.707	5.208	5.959
7	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	1.397	1.86	2.306	2.896	3.355	4.501	5.041
9	1.383	1.833	2.262	2.821	3.25	4.297	4.781
10	1.372	1.812	2.228	2.764	3.169	4.144	4.587
11	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	1.356	1.782	2.179	2.681	3.055	3.93	4.318
13	1.35	1.771	2.16	2.65	3.012	3.852	4.221
14	1.345	1.761	2.145	2.624	2.977	3.787	4.14
15	1.341	1.753	2.131	2.602	2.947	3.733	4.073
120	1.289	1.658	1.98	2.358	2.617	3.16	3.373
1000	1.282	1.646	1.962	2.33	2.581	3.098	3.3
Inf	1.282	1.645	1.96	2.326	2.576	3.091	3.291

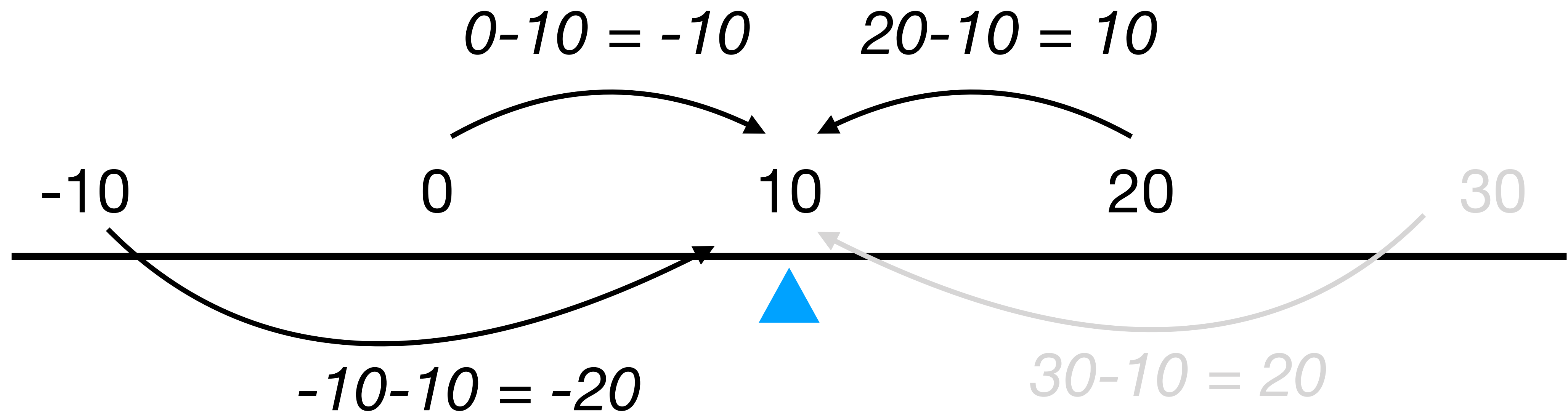
# N-1: degrees of freedom explained



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

**Sum of deviations is always zero!**

# N-1: degrees of freedom explained



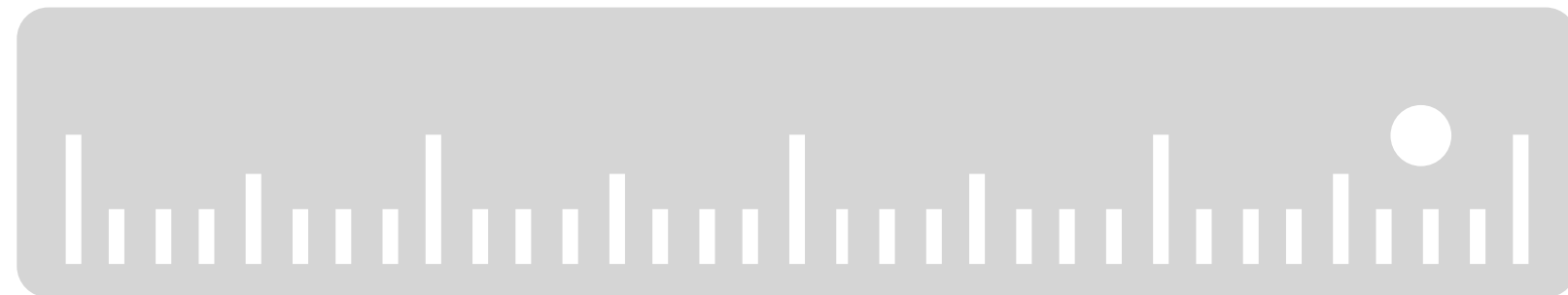
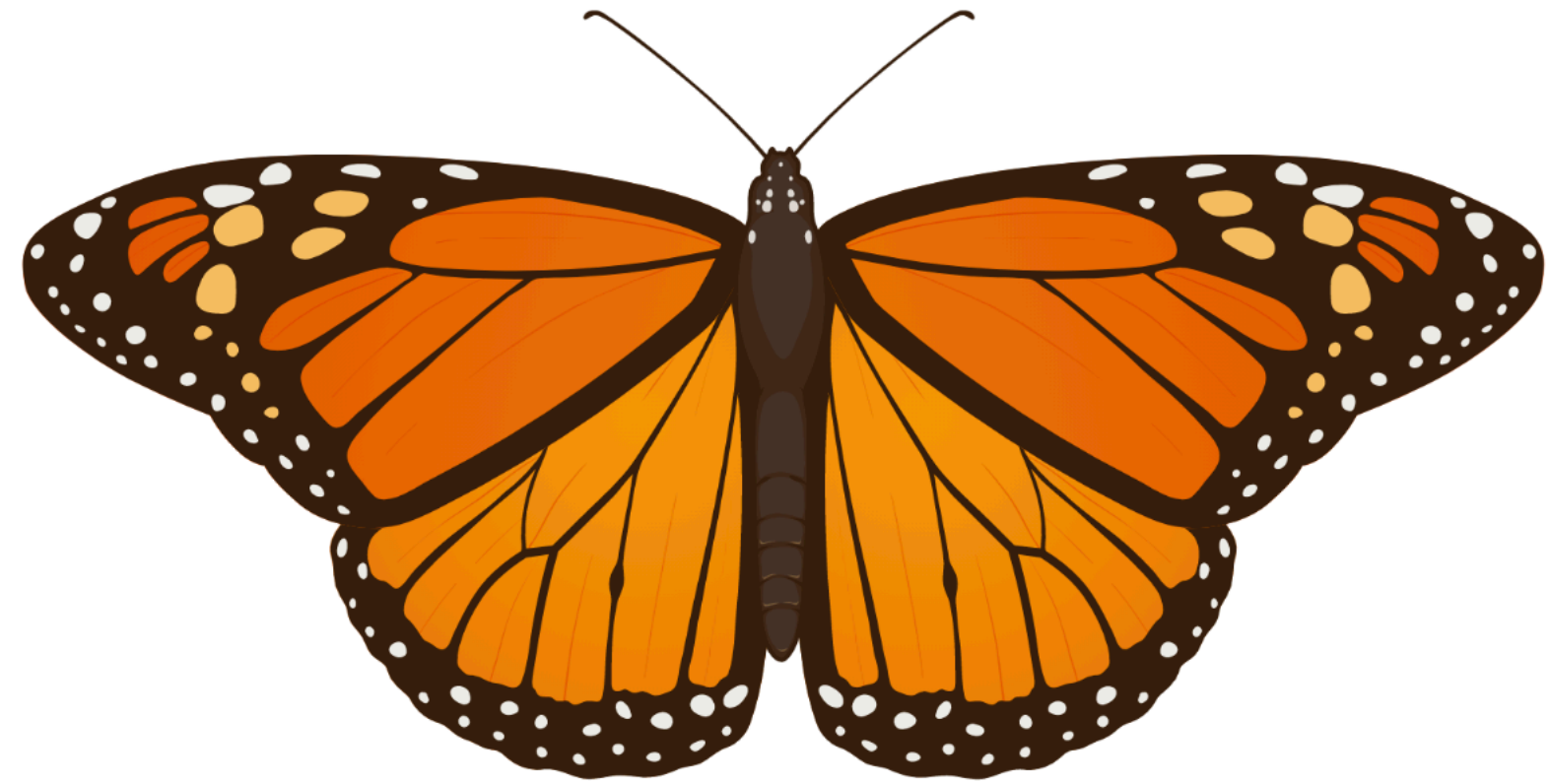
$$(-20) + (-10) + (10) + \mathbf{X} = 0$$

Sum of deviations is always zero!

Has to be  
+20 from  
mean...



# Calculating the confidence interval: butterflies



$$n = 14$$

$$df = 13$$

$$\bar{y} = 32.81 \text{ cm}^2$$

$$s = 2.48 \text{ cm}^2$$

$$\bar{y} \pm t_{0.025} \frac{s}{\sqrt{n}}$$

	P						
one-tail	0.1	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	0.2	0.1	0.05	0.02	0.01	0.002	0.001
DF							
1	3.078	6.314	12.706	31.821	63.656	318.289	636.578
2	1.886	2.92	4.303	6.965	9.925	22.328	31.6
3	1.638	2.353	3.182	4.541	5.841	10.214	12.924
4	1.533	2.132	2.776	3.747	4.604	7.173	8.61
5	1.476	2.015	2.571	3.365	4.032	5.894	6.869
6	1.44	1.943	2.447	3.143	3.707	5.208	5.959
7	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	1.397	1.86	2.306	2.896	3.355	4.501	5.041
9	1.383	1.833	2.262	2.821	3.25	4.297	4.781
10	1.372	1.812	2.228	2.764	3.169	4.144	4.587
11	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	1.356	1.782	2.179	2.681	3.055	3.93	4.318
13	1.35	1.771	2.16	2.65	3.012	3.852	4.221
14	1.345	1.761	2.145	2.624	2.977	3.787	4.14
15	1.341	1.753	2.131	2.602	2.947	3.733	4.073
120	1.289	1.658	1.98	2.358	2.617	3.16	3.373
1000	1.282	1.646	1.962	2.33	2.581	3.098	3.3
Inf	1.282	1.645	1.96	2.326	2.576	3.091	3.291

# Calculating the confidence interval: butterflies



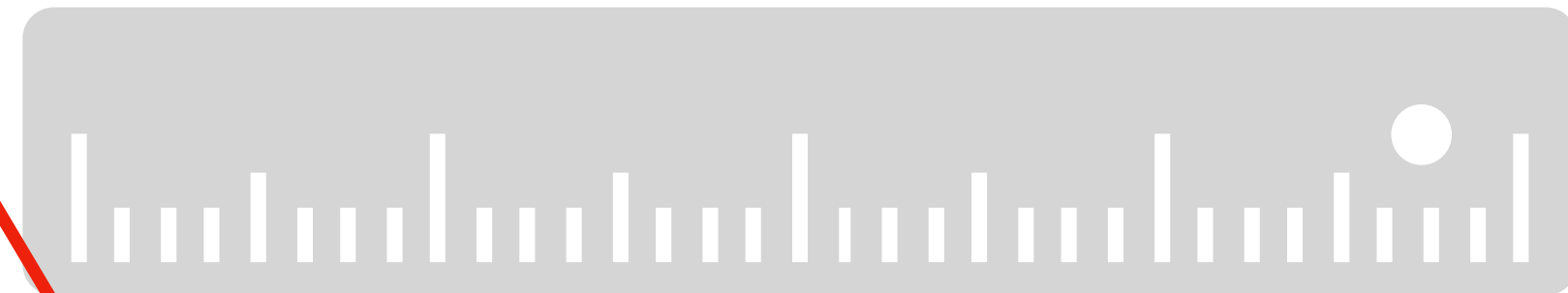
$$n = 14$$

$$df = 13$$

$$\bar{y} = 32.81 \text{ cm}^2$$

$$s = 2.48 \text{ cm}^2$$

**Critical  
value**



**qt (p, df)**

qt (0.975, 13)

**95%**

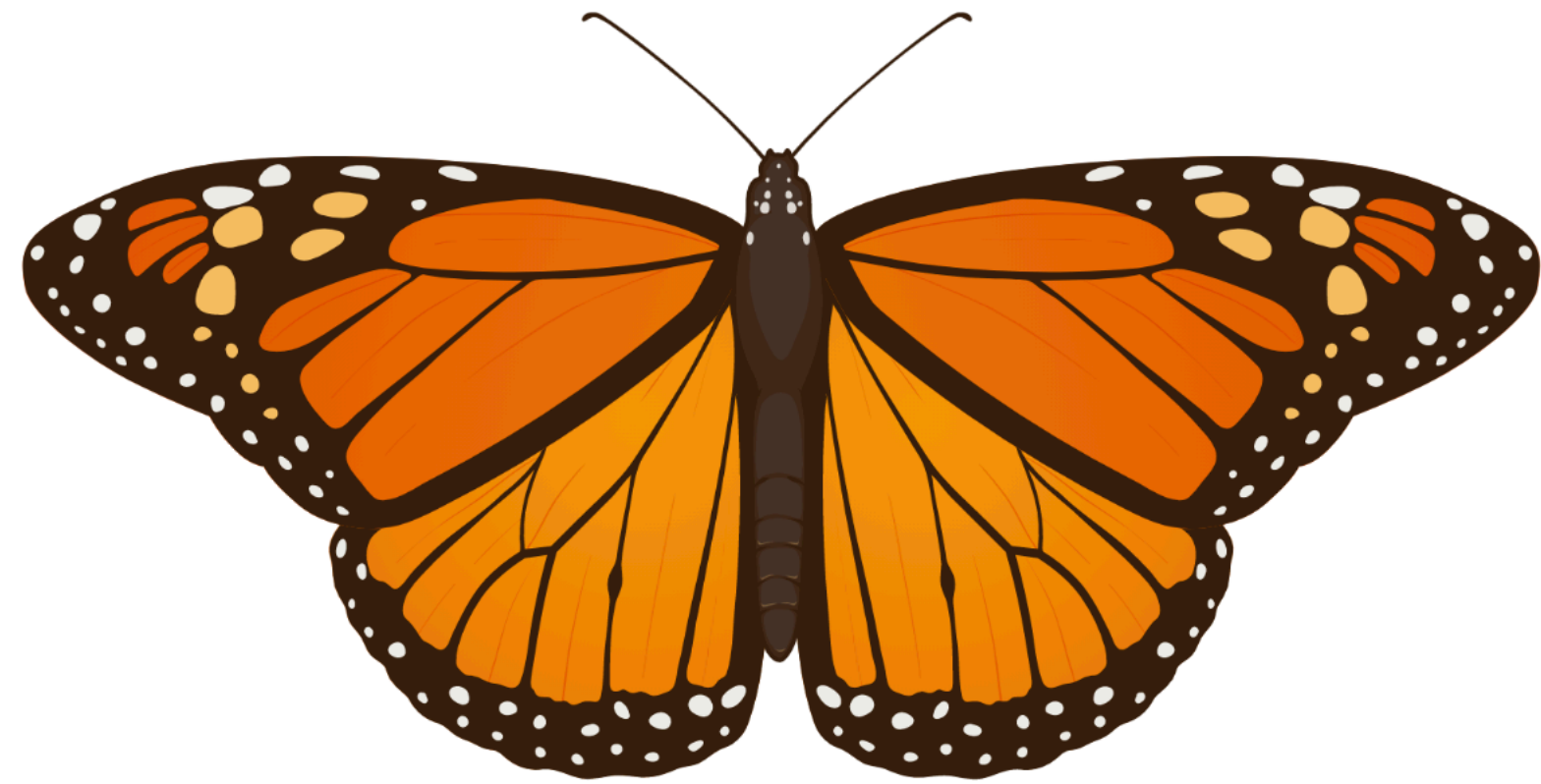
$$\bar{y} \pm t_{0.025} \frac{s}{\sqrt{n}}$$

$$32.81 \pm 2.16 \frac{2.48}{\sqrt{14}}$$

$$32.81 \pm 1.43 \quad (31.4, 34.2)$$



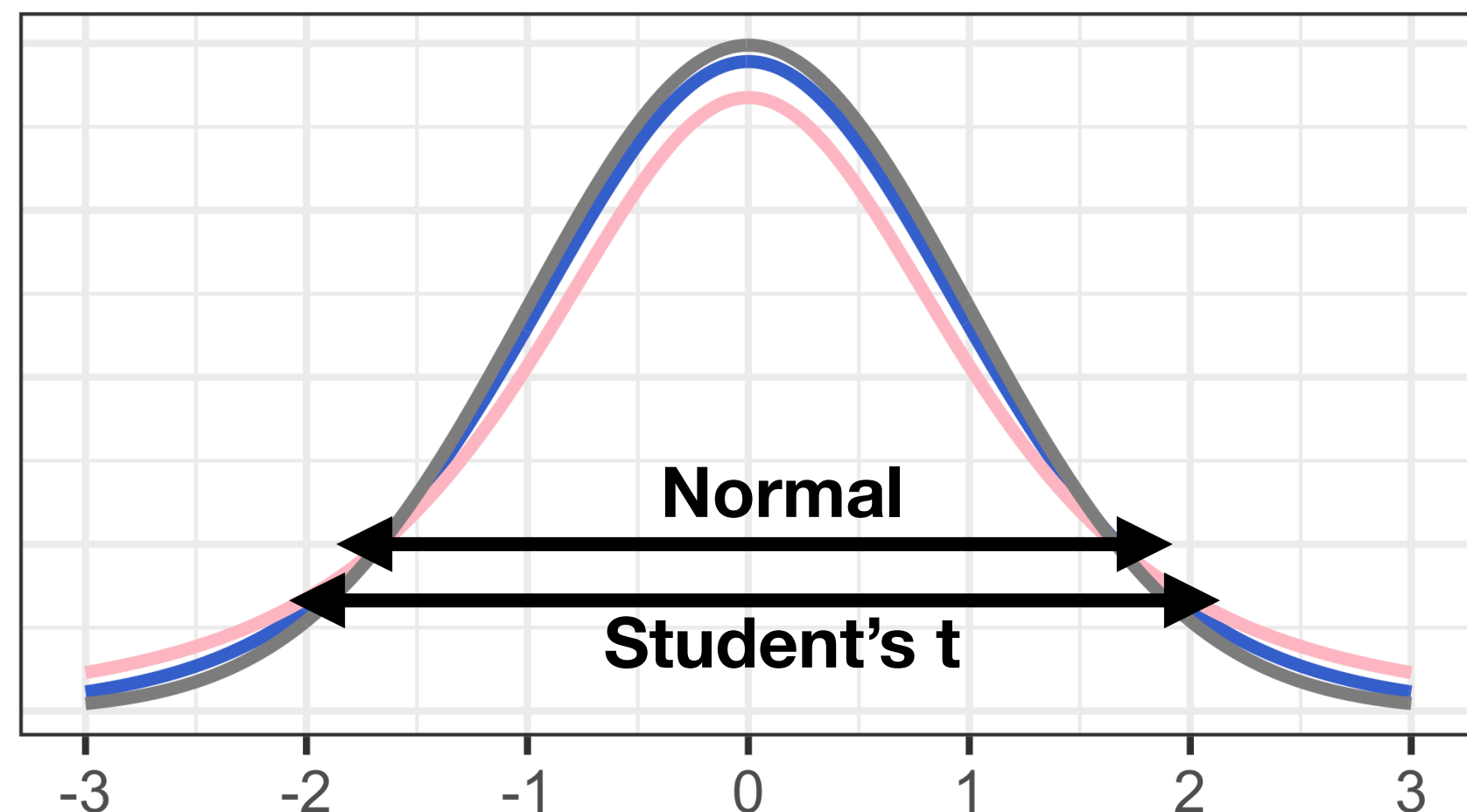
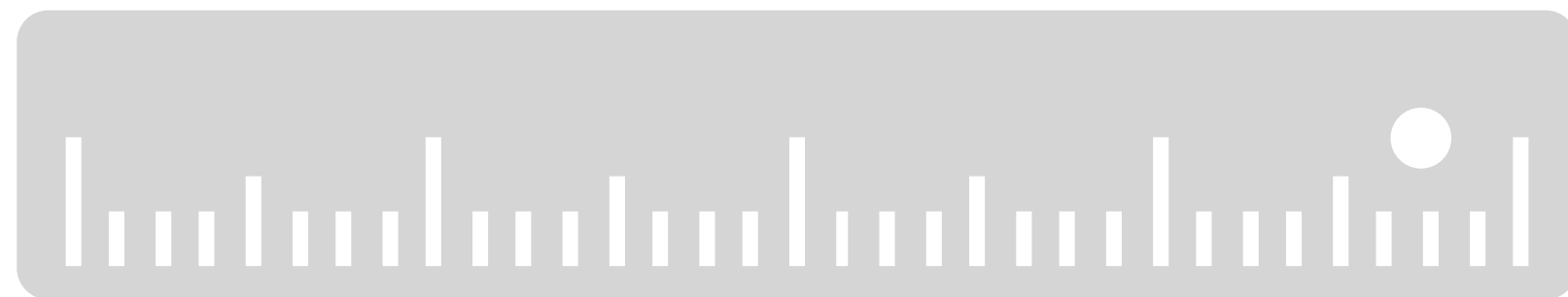
# Note: why use Student's $t$ distribution?



$$\begin{array}{c} n = 14 \\ \longrightarrow \\ df = 13 \end{array}$$

$$\bar{y} = 32.81 \text{ cm}^2$$

$$s = 2.48 \text{ cm}^2$$



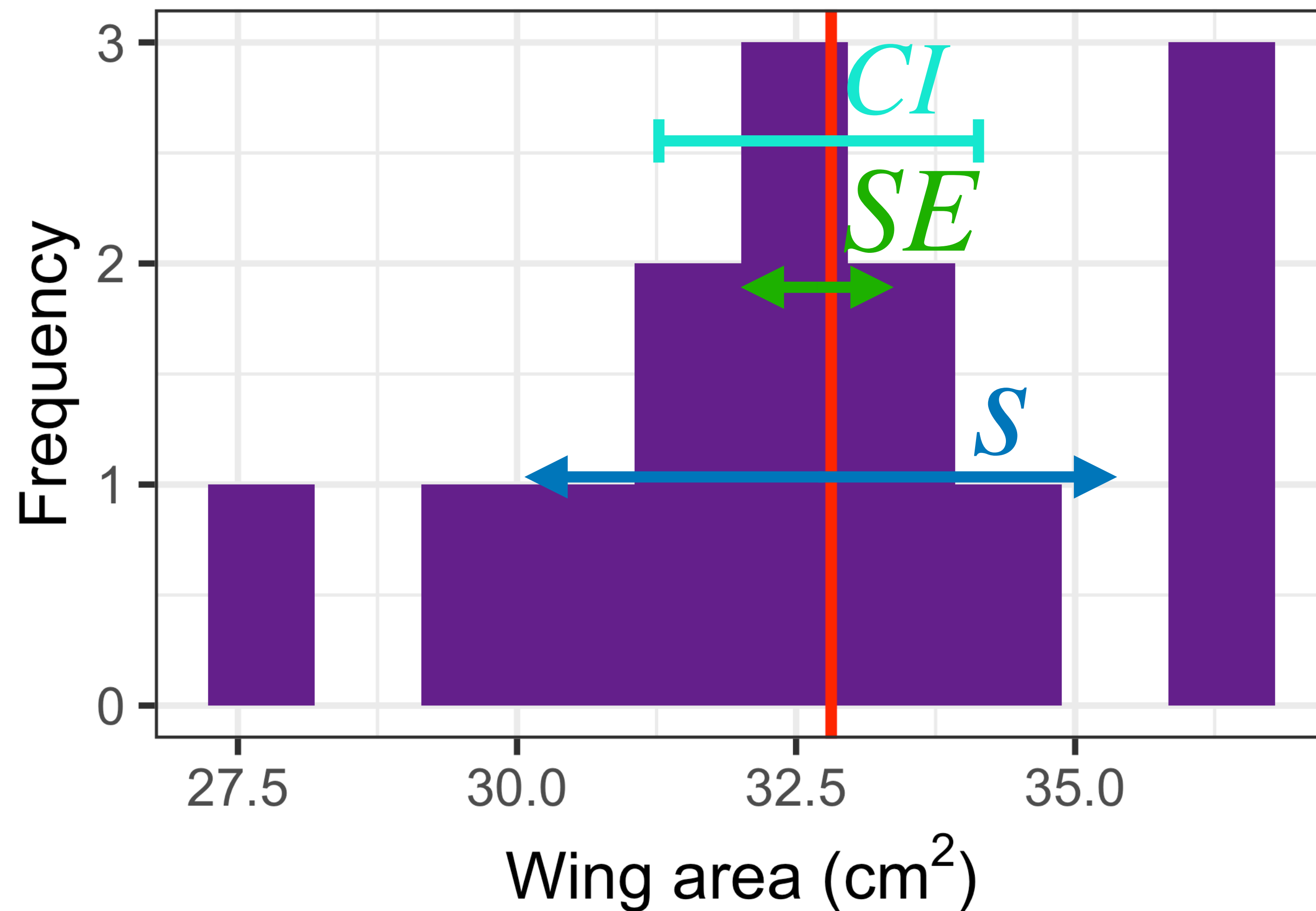
95%

$$\bar{y} \pm z_{0.025} \frac{s}{\sqrt{n}}$$

$$32.81 \pm 1.96 \frac{2.48}{\sqrt{14}}$$

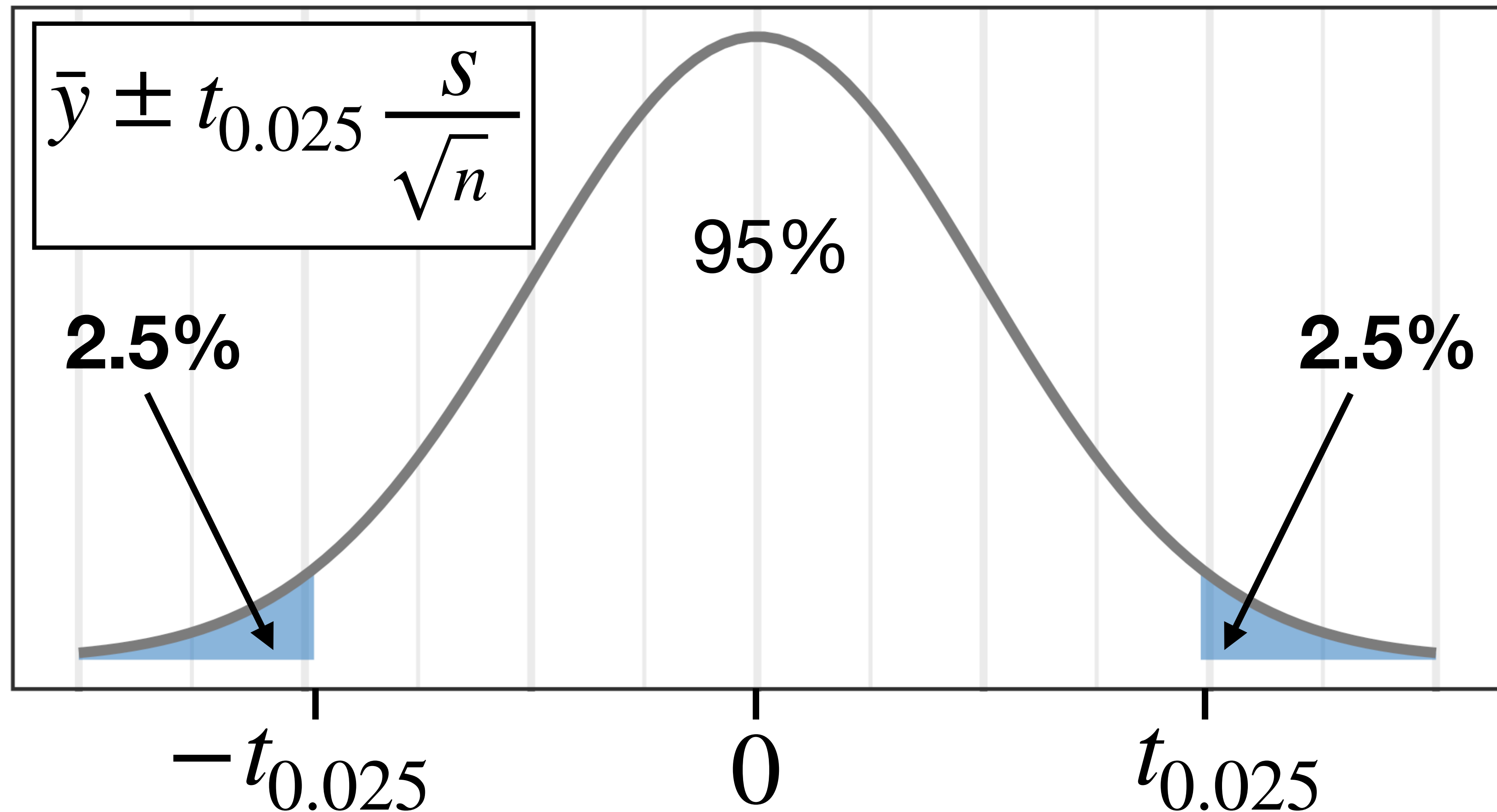
$$32.81 \pm 1.29 \quad (31.51, 34.11)$$

# Calculating the confidence interval: butterflies

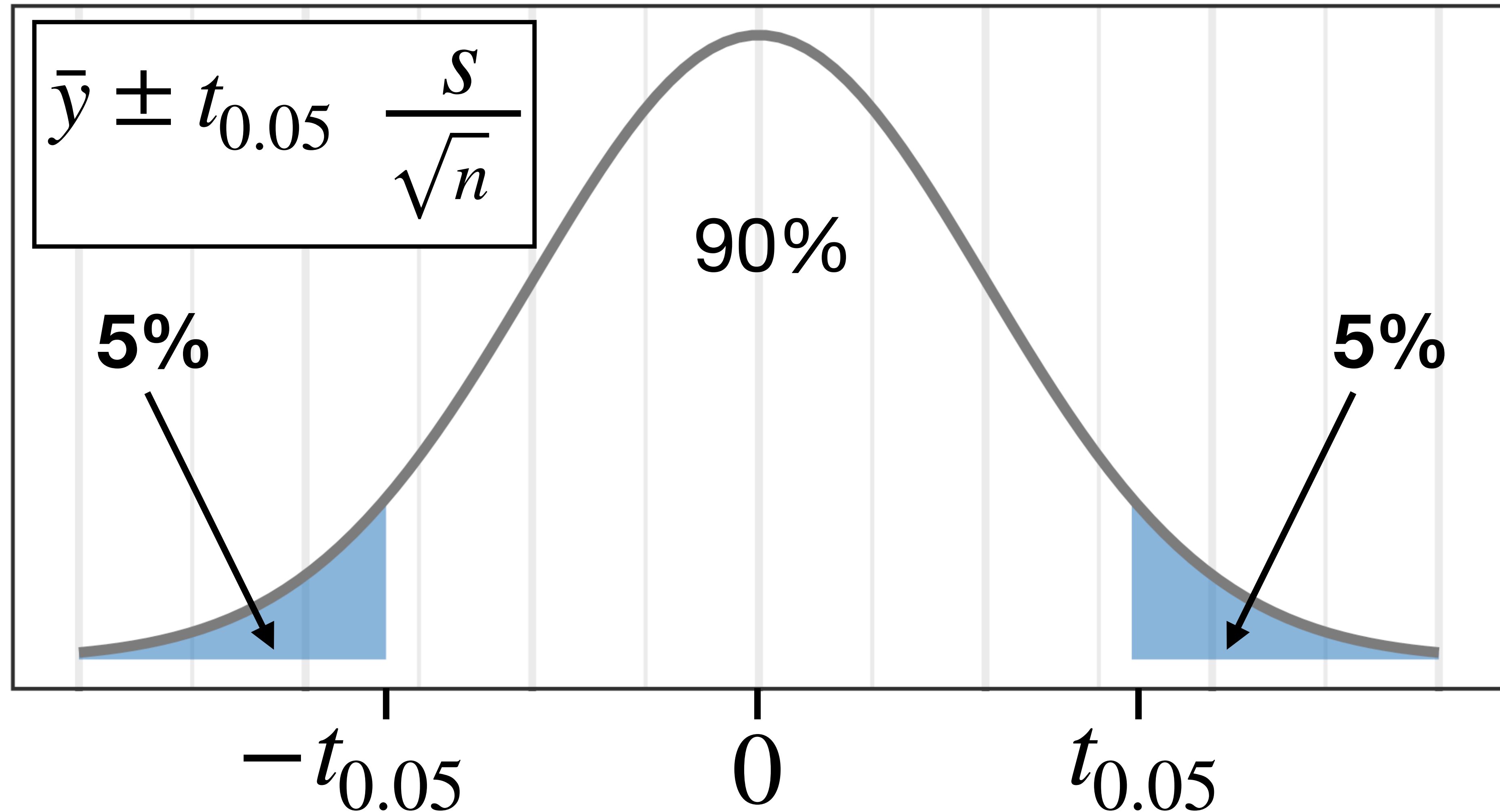


$32.81 \pm 1.43$

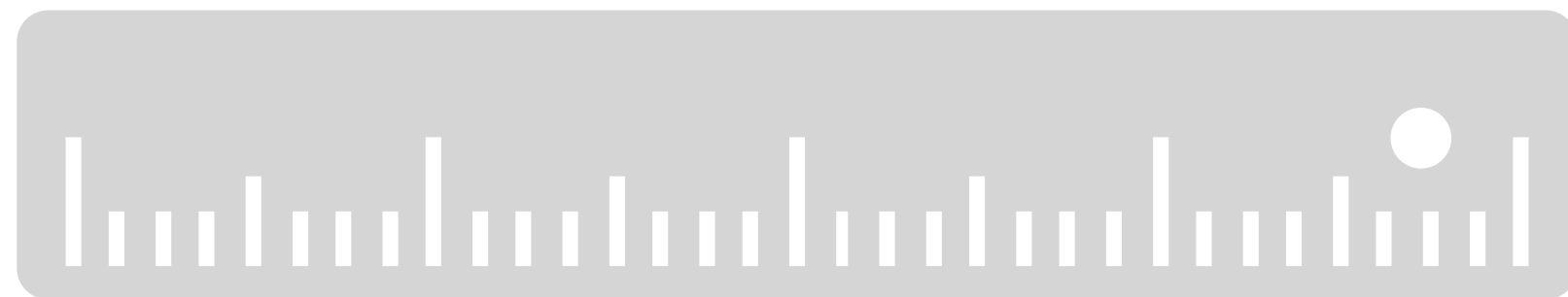
# Critical value and Student's $t$ distribution



# Critical value and Student's $t$ distribution



# Calculating the confidence interval: butterflies



$$n = 14$$

$$df = 13$$

$$\bar{y} = 32.81 \text{ cm}^2$$

$$s = 2.48 \text{ cm}^2$$

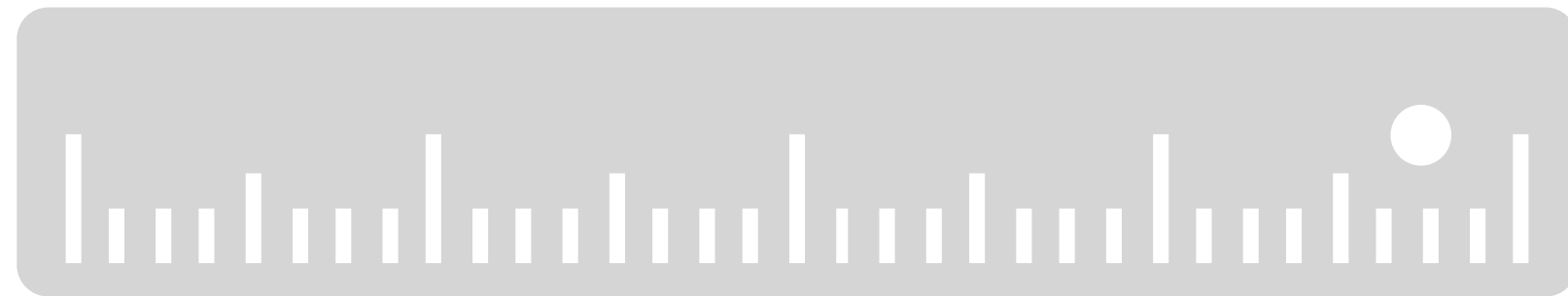
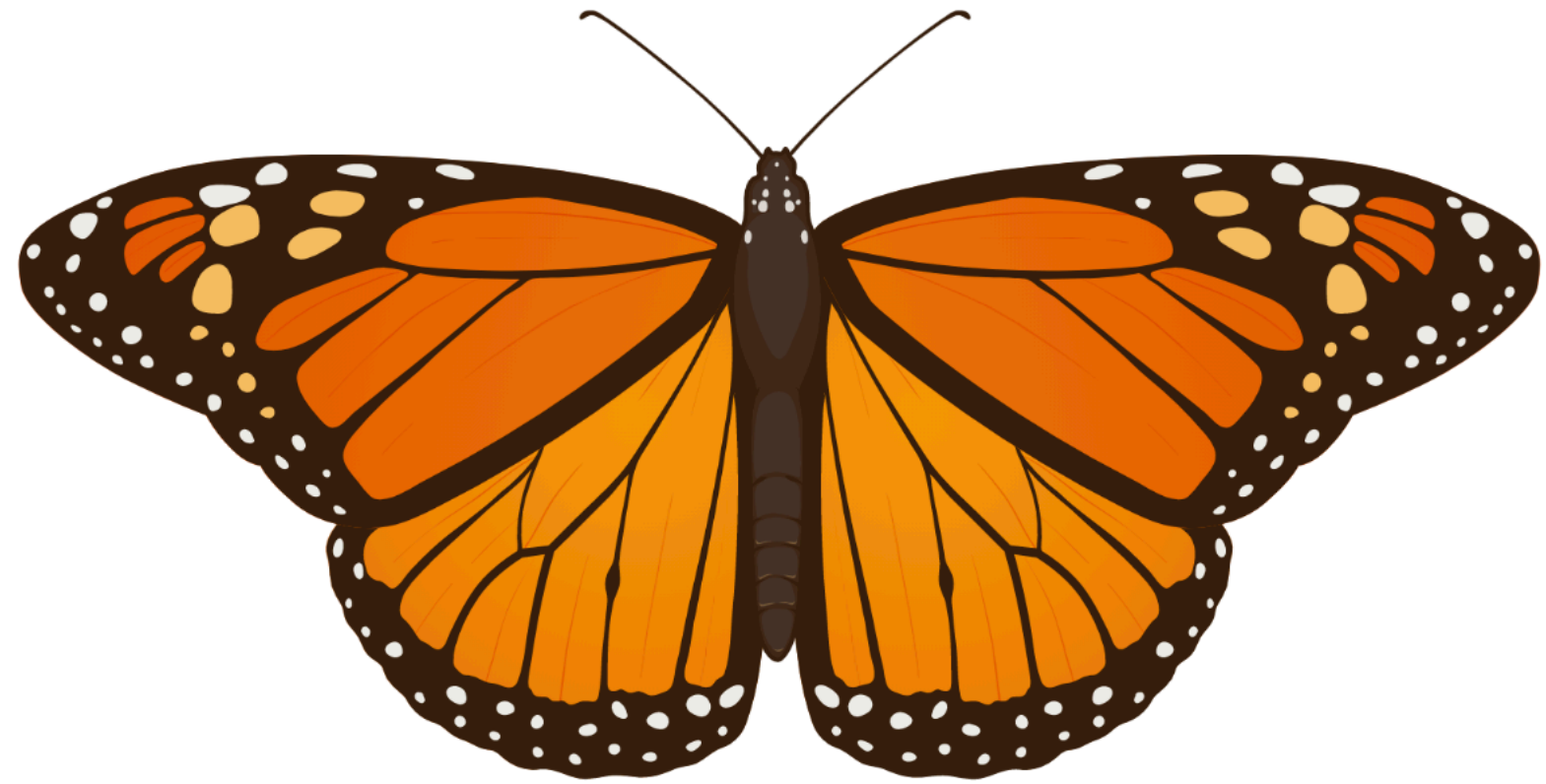
90%

$$\bar{y} \pm t_{0.05} \frac{s}{\sqrt{n}}$$



	P						
one-tail	0.1	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	0.2	0.1	0.05	0.02	0.01	0.002	0.001
DF							
1	3.078	6.314	12.706	31.821	63.656	318.289	636.578
2	1.886	2.92	4.303	6.965	9.925	22.328	31.6
3	1.638	2.353	3.182	4.541	5.841	10.214	12.924
4	1.533	2.132	2.776	3.747	4.604	7.173	8.61
5	1.476	2.015	2.571	3.365	4.032	5.894	6.869
6	1.44	1.943	2.447	3.143	3.707	5.208	5.959
7	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	1.397	1.86	2.306	2.896	3.355	4.501	5.041
9	1.383	1.833	2.262	2.821	3.25	4.297	4.781
10	1.372	1.812	2.228	2.764	3.169	4.144	4.587
11	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	1.356	1.782	2.179	2.681	3.055	3.93	4.318
13	1.35	1.771	2.16	2.65	3.012	3.852	4.221
14	1.345	1.761	2.145	2.624	2.977	3.787	4.14
15	1.341	1.753	2.131	2.602	2.947	3.733	4.073
120	1.289	1.658	1.98	2.358	2.617	3.16	3.373
1000	1.282	1.646	1.962	2.33	2.581	3.098	3.3
Inf	1.282	1.645	1.96	2.326	2.576	3.091	3.291

# Calculating the confidence interval: butterflies



**The higher the confidence level, the wider the confidence interval**

$$n = 14$$

$$df = 13$$

$$\bar{y} = 32.81 \text{ cm}^2$$

$$s = 2.48 \text{ cm}^2$$

**90%**

$$\bar{y} \pm t_{0.05} \frac{s}{\sqrt{n}}$$

$$32.81 \pm 1.77 \frac{2.48}{\sqrt{14}}$$

$$32.81 \pm 1.17 \quad (31.6, 34.0)$$



# Example

A pharmacologist measured the concentration of dopamine in the brains of eight rats. The mean concentration was 1,269 ng/gm and the standard deviation was 145 ng/gm. **Construct a 95% confidence interval for the population mean.**

$$SE = \frac{s}{\sqrt{n}}$$

$$SE = \frac{145}{\sqrt{8}} = 51.2$$

$$\bar{y} \pm t_{0.025} \frac{s}{\sqrt{n}} \quad (\text{df} = n - 1 = 7)$$

$$1269 \pm (t_{0.025})(51.2)$$

	P						
one-tail	0.1	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	0.2	0.1	0.05	0.02	0.01	0.002	0.001
DF							
1	3.078	6.314	12.706	31.821	63.656	318.289	636.578
2	1.886	2.92	4.303	6.965	9.925	22.328	31.6
3	1.638	2.353	3.182	4.541	5.841	10.214	12.924
4	1.533	2.132	2.776	3.747	4.604	7.173	8.61
5	1.476	2.015	2.571	3.365	4.032	5.894	6.869
6	1.44	1.943	2.447	3.143	3.707	5.208	5.959
7	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	1.397	1.86	2.306	2.896	3.355	4.501	5.041
9	1.383	1.833	2.262	2.821	3.25	4.297	4.781
10	1.371	1.812	2.228	2.764	3.183	4.144	4.587
11	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	1.356	1.782	2.177	2.675	3.032	3.93	4.318
13	1.35	1.771	2.156	2.635	2.962	3.852	4.221
14	1.345	1.761	2.145	2.624	2.977	3.787	4.14
15	1.341	1.753	2.131	2.602	2.947	3.733	4.073
20	1.329	1.729	2.086	2.528	2.878	3.619	3.93
30	1.315	1.704	2.042	2.457	2.819	3.527	3.815
40	1.308	1.69	2.019	2.415	2.781	3.465	3.745
60	1.299	1.674	1.995	2.368	2.733	3.406	3.686
120	1.289	1.658	1.98	2.358	2.617	3.16	3.373
1000	1.282	1.646	1.962	2.33	2.581	3.098	3.3
Inf	1.282	1.645	1.96	2.326	2.576	3.091	3.291

qt(p = 0.025, df = 7, lower.tail = F)

qt(p = 0.975, df = 7, lower.tail = T)

# Example

A pharmacologist measured the concentration of dopamine in the brains of eight rats. The mean concentration was 1,269 ng/gm and the standard deviation was 145 ng/gm. **Construct a 95% confidence interval for the population mean.**

$$SE = \frac{s}{\sqrt{n}}$$

$$SE = \frac{145}{\sqrt{8}} = 51.2$$

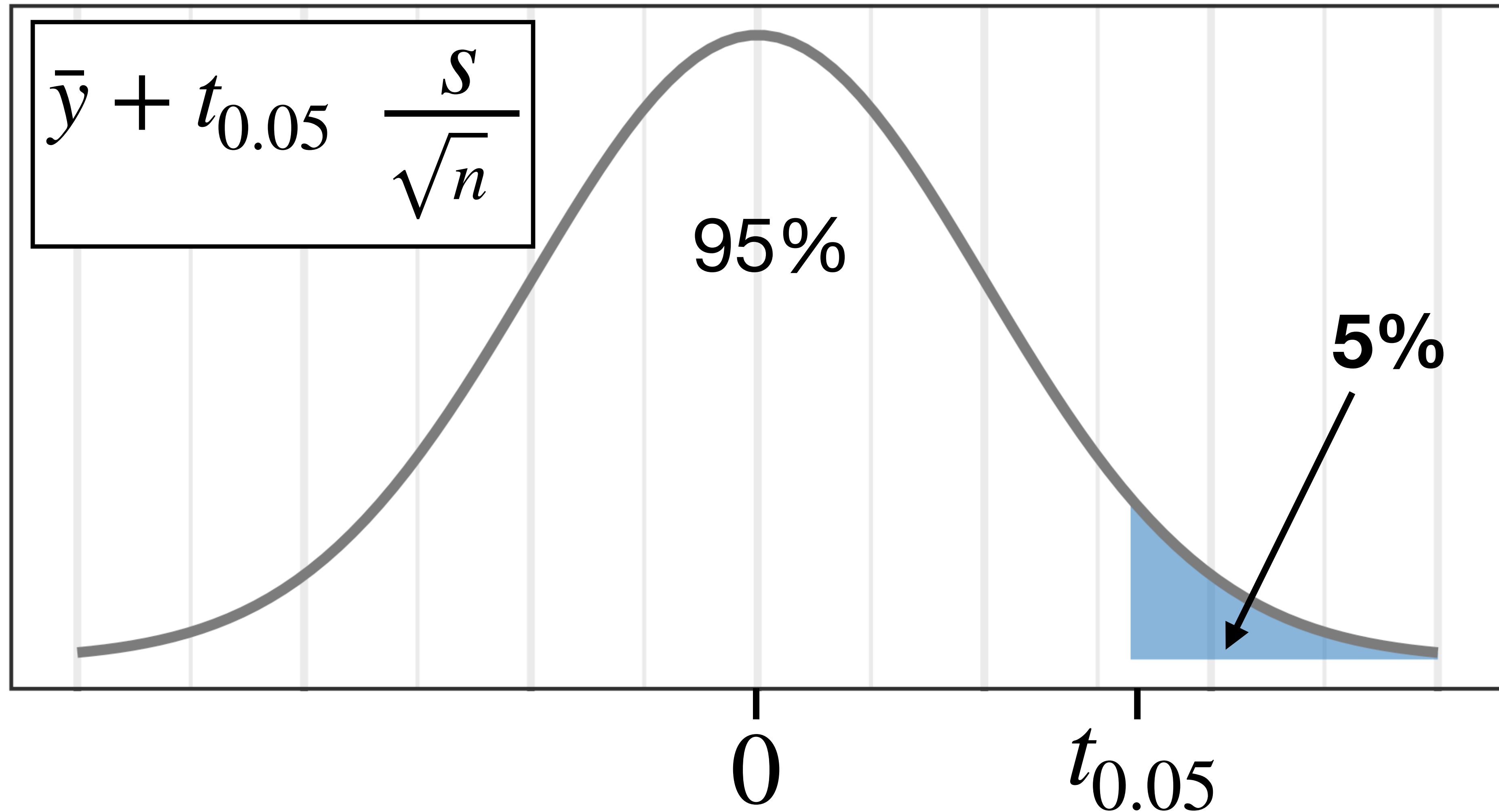
$$\bar{y} \pm t_{0.025} \frac{s}{\sqrt{n}} \quad (\text{df} = n - 1 = 7)$$

$$1269 \pm (2.365)(51.2)$$

$$1269 \pm 121.1$$

$$(1147.9, 1390.1)$$

# One-sided confidence intervals

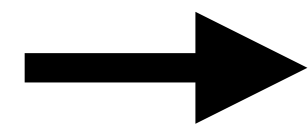


# One-sided confidence intervals

A pharmacologist measured the concentration of dopamine in the brains of eight rats. The mean concentration was 1,269 ng/gm and the standard deviation was 145 ng/gm. **Construct a one-sided, upper-bound 95% confidence interval for the population mean.**

$$\bar{y} + t_{0.05} \frac{S}{\sqrt{n}} \quad (\text{df} = n - 1 = 7)$$





	P						
one-tail	0.1	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	0.2	0.1	0.05	0.02	0.01	0.002	0.001
DF							
1	3.078	6.314	12.706	31.821	63.656	318.289	636.578
2	1.886	2.92	4.303	6.965	9.925	22.328	31.6
3	1.638	2.353	3.182	4.541	5.841	10.214	12.924
4	1.533	2.132	2.776	3.747	4.604	7.173	8.61
5	1.476	2.015	2.571	3.365	4.032	5.894	6.869
6	1.44	1.943	2.447	3.143	3.707	5.208	5.959
7	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	1.397	1.86	2.306	2.896	3.355	4.501	5.041
9	1.383	1.833	2.262	2.821	3.25	4.297	4.781
10	1.371	1.812	2.228	2.764	3.183	4.144	4.587
11	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	1.357	1.782	2.179	2.675	3.037	3.93	4.318
13	1.353	1.771	2.16	2.65	3.012	3.852	4.221
14	1.345	1.761	2.145	2.624	2.977	3.787	4.14
15	1.341	1.753	2.131	2.602	2.947	3.733	4.073
20	1.328	1.729	2.086	2.539	2.876	3.579	3.858
30	1.315	1.704	2.045	2.479	2.812	3.450	3.686
40	1.306	1.688	2.015	2.437	2.764	3.347	3.579
50	1.299	1.676	1.992	2.407	2.729	3.267	3.494
60	1.294	1.666	1.975	2.385	2.701	3.203	3.431
70	1.29	1.658	1.962	2.368	2.678	3.147	3.373
80	1.287	1.652	1.952	2.354	2.658	3.097	3.327
90	1.285	1.647	1.944	2.342	2.641	3.054	3.287
100	1.282	1.645	1.938	2.332	2.627	3.019	3.254
120	1.282	1.645	1.936	2.328	2.622	3.009	3.241
1000	1.282	1.645	1.936	2.326	2.576	3.091	3.291
Inf	1.282	1.645	1.936	2.326	2.576	3.091	3.291

`qt(p = 0.05, df = 7, lower.tail = F)`

`qt(p = 0.95, df = 7, lower.tail = T)`

# One-sided confidence intervals

A pharmacologist measured the concentration of dopamine in the brains of eight rats. The mean concentration was 1,269 ng/gm and the standard deviation was 145 ng/gm. **Construct a one-sided, upper-bound 95% confidence interval for the population mean.**

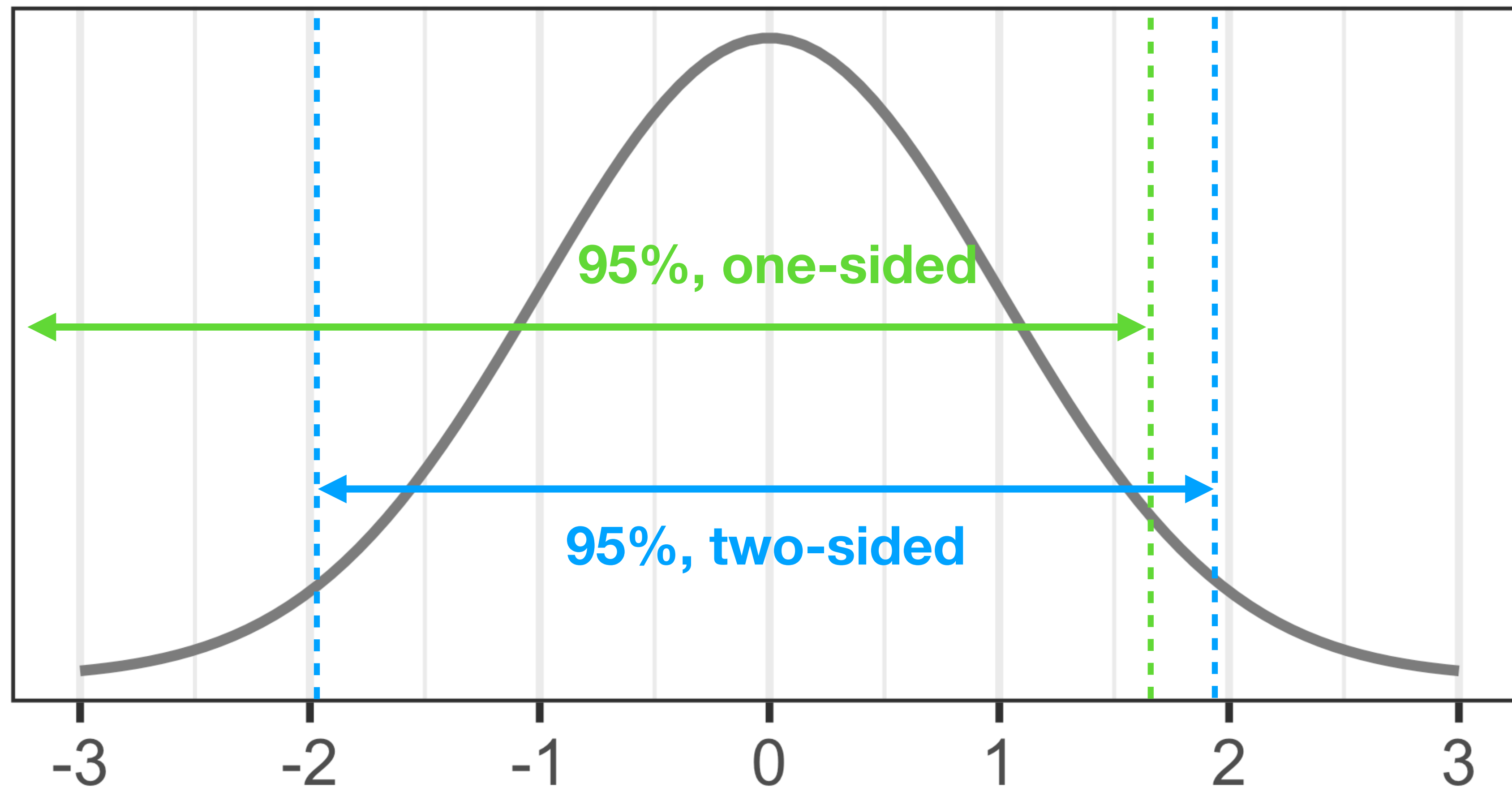
$$\bar{y} + t_{0.05} \frac{s}{\sqrt{n}} \quad (\text{df} = n - 1 = 7)$$

$$1269 + (1.895)(51.2)$$

$$(\infty, 1366.0)$$



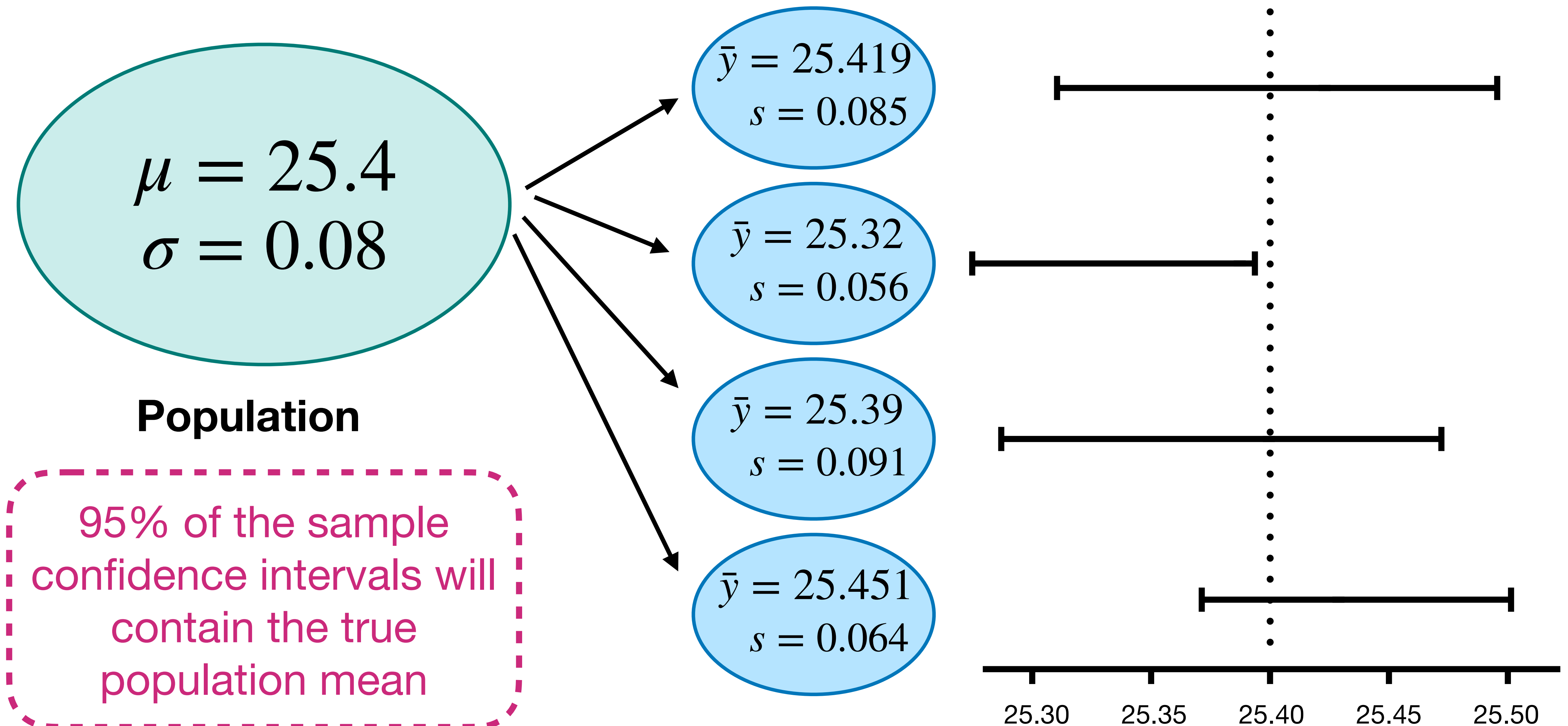
# One-sided confidence intervals



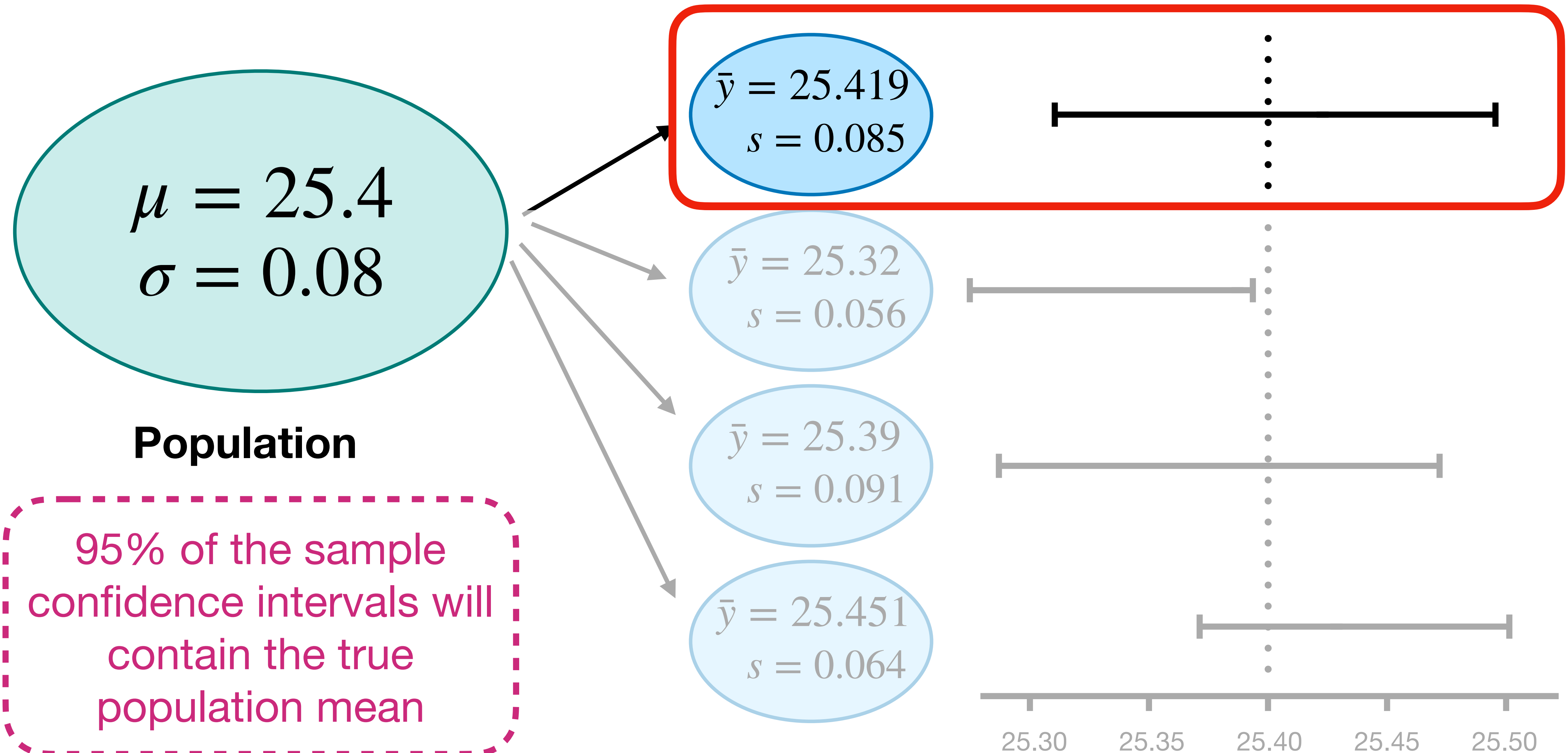
$(-\infty, 1366.0)$

$(1147.9, 1390.1)$

# Confidence intervals and randomness



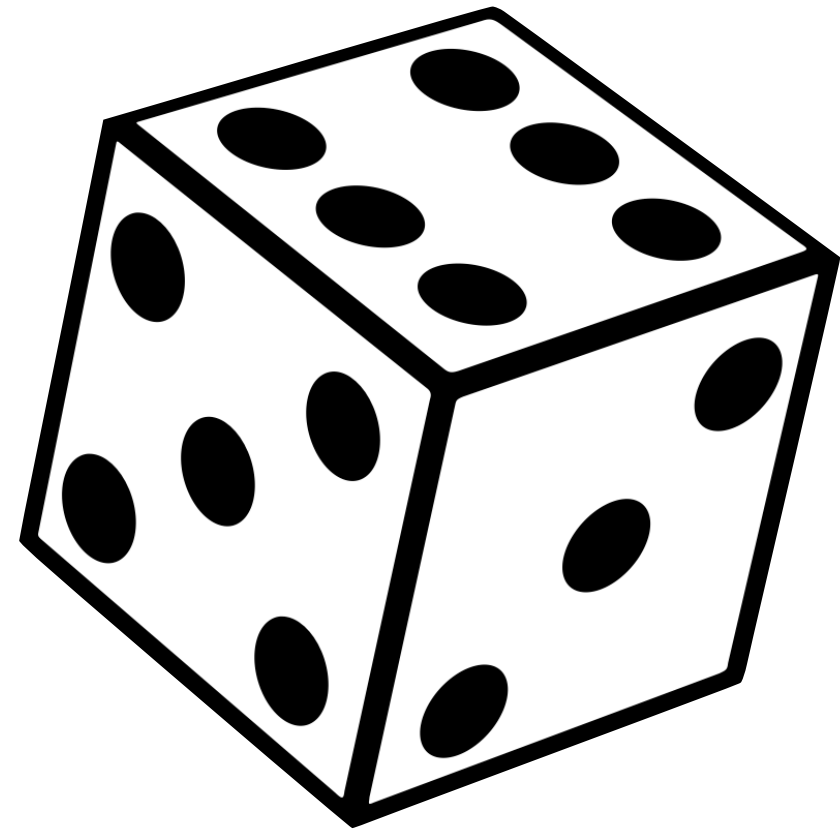
# Confidence intervals and randomness



# Confidence intervals and randomness

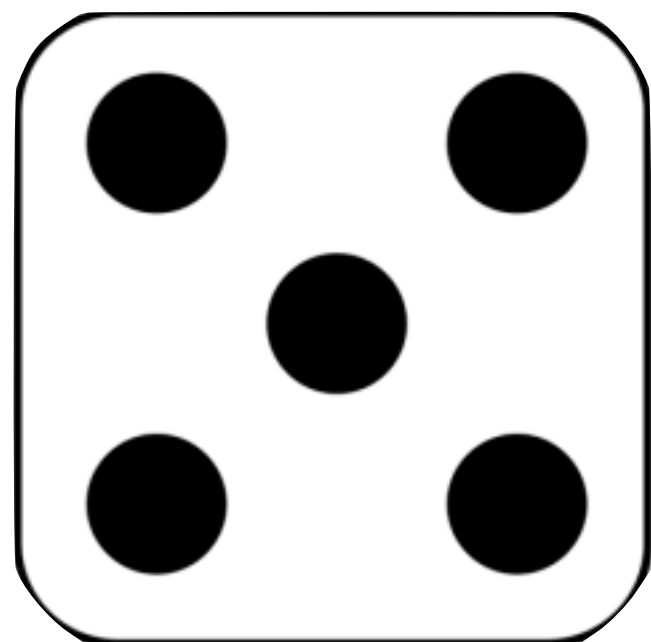
- **Larger samples** produce **narrower** confidence intervals
  - *Because SE is smaller (divide by square root of  $n$ )*
- A **confidence interval** can be interpreted as a probability... with caution!

# Confidence intervals and randomness



$$Pr[Y = 2] = \frac{1}{6}$$

Pr{a sample will give us a CI that contains the true mean} = 0.95



$$Y = 5$$

$$Pr[5 = 2] \neq \frac{1}{6}$$

Pr{the true mean is within our CI} = 0.95



$$Pr\{31 < \mu < 34\} \neq 0.95$$

# Confidence intervals and randomness

- **Larger samples** produce **narrower** confidence intervals
  - *Because SE is smaller (divide by square root of  $n$ )*
- A **confidence interval** can be interpreted as a probability... with caution!
  - $\Pr\{\text{a sample will give us a CI that contains the true mean}\} = 0.95$  ✓
  - $\Pr\{\text{the true mean is within our CI}\} = 0.95$  ✗
  - *An individual statement can be TRUE or FALSE, but if you create numerous statements, one statement will be TRUE 95% of the time*
- “We are 95% confident that the true mean is between X and X”

A pharmacologist measured the concentration of dopamine in the brains of eight rats. The mean concentration was 1,269 ng/gm and the standard deviation was 145 ng/gm. **Construct a 95% confidence interval for the population mean.**

$$SE = \frac{s}{\sqrt{n}}$$

$$SE = \frac{145}{\sqrt{8}} = 51.2$$

$$\bar{y} \pm t_{0.025} \frac{s}{\sqrt{n}} \quad (\text{df} = n - 1 = 7)$$

$$1269 \pm (2.365)(51.2)$$

$$1269 \pm 121.1$$

$$(1147.9, 1390.1)$$

**We are 95% confident that the mean concentration of dopamine of all rats is between 1,147.9 and 1,390.1 ng/gm**



# Planning a study with sufficient precision

**There is no substitute for bad data.**

- **(1) Population variability of observed variable**
  - Sometimes you cannot (or should not) reduce this, you want a random sample of the entire population
  - However, reducing environmental variation and fixing the variables you can will lead to the cleanest data

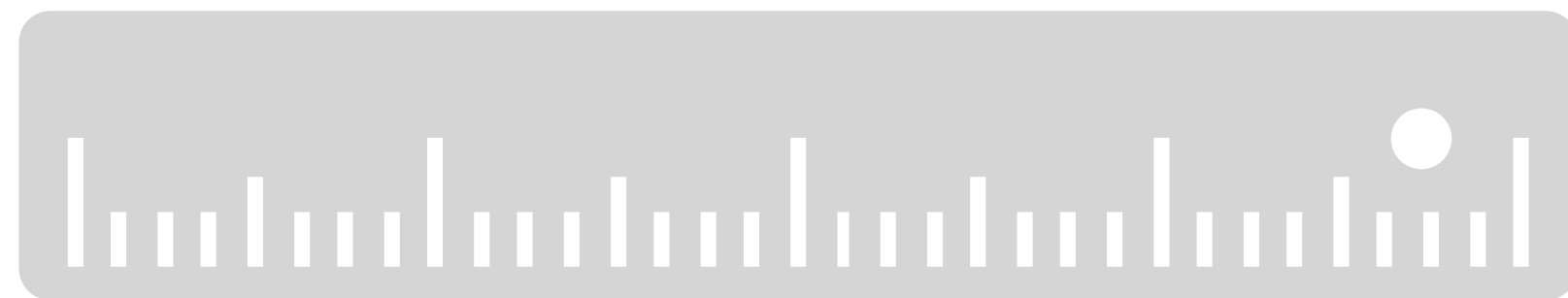


# Planning a study with sufficient precision

**There is no substitute for bad data.**

- **(1) Population variability of observed variable**
  - Sometimes you cannot (or should not) reduce this, you want a random sample of the entire population
  - However, reducing environmental variation and fixing the variables you can will lead to the cleanest data
- **(2) Sample size**
  - You can try to guess-calculate how large of a sample you need to prove/disprove your hypothesis
  - $SE = s/\sqrt{n}$

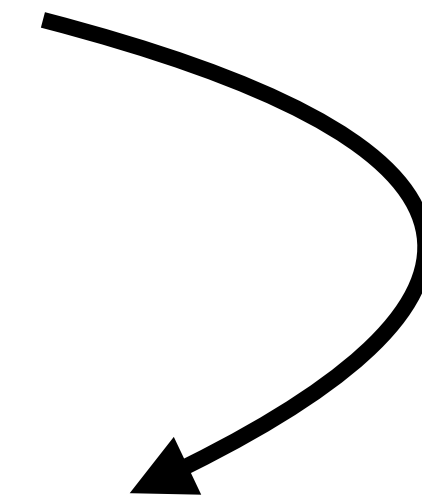
# Planning a study with sufficient precision



$$\xrightarrow{n = 14} \begin{aligned} \bar{y} &= 32.81 \text{ cm}^2 \\ s &= 2.48 \text{ cm}^2 \end{aligned}$$

$$SE = \frac{2.48}{\sqrt{n}} \leq 0.4$$

$$n \geq 38.4$$



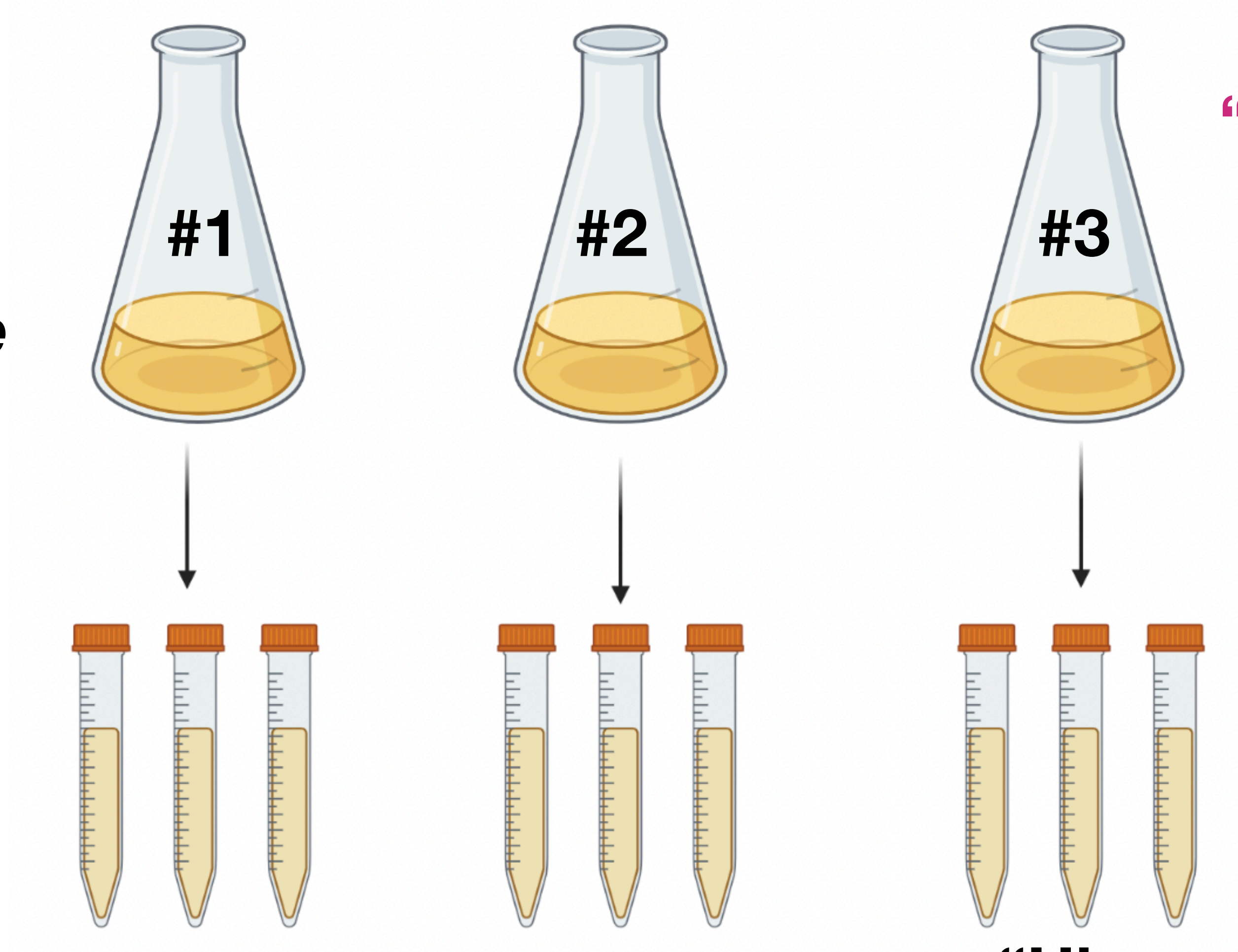
# Assumptions for estimating confidence intervals

- **Conditions on the design of the study:**
  - (1) Data is a random sample from a large population
  - (2) Observations in the sample must be independent of each other

# What does it mean for samples to be independent of each other?

*Still a good experimental design! Means > individual sample*

**$n = 3?$**



**“Biological replicates”**

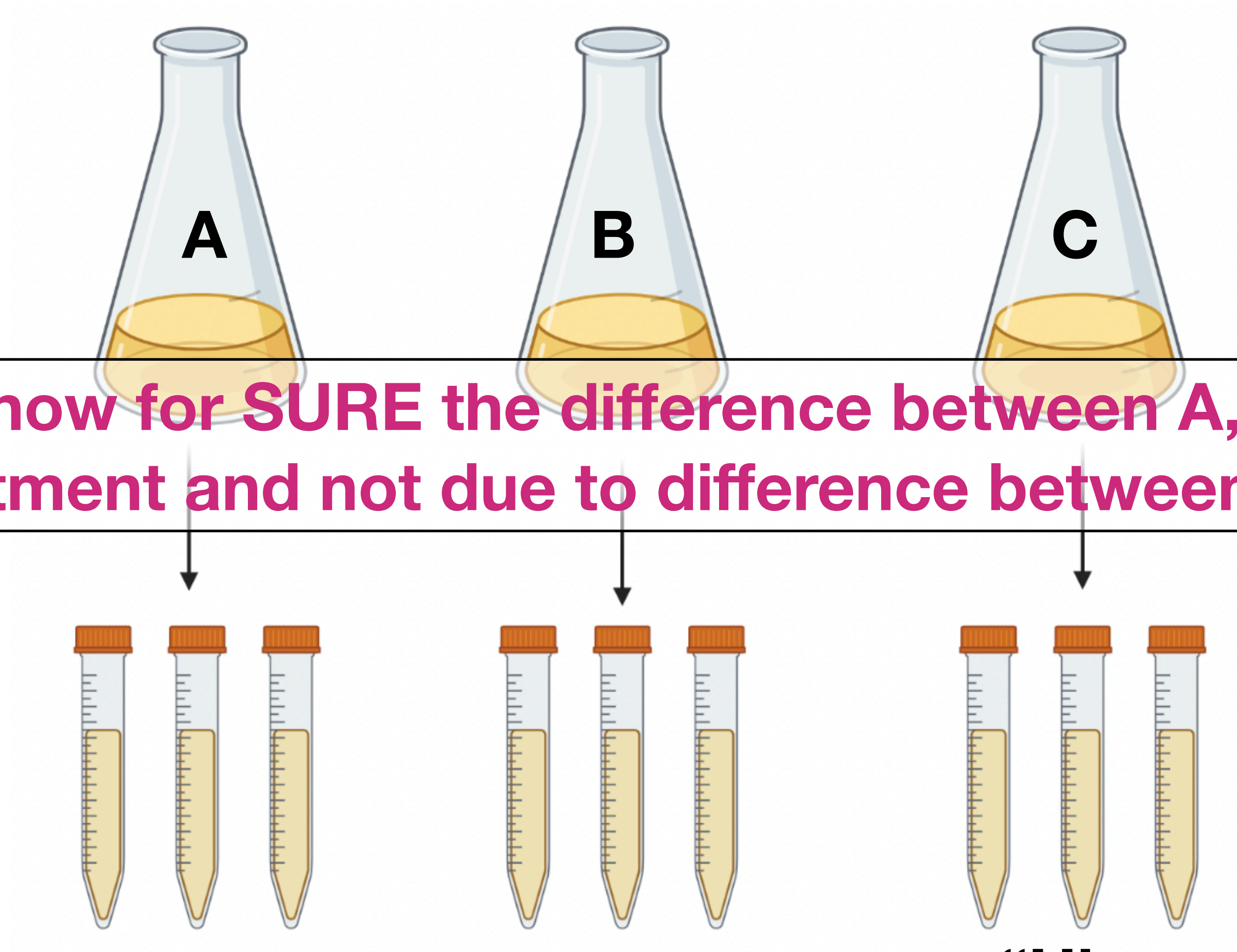
**$n = 9?$**

**“Technical replicates”**

**“Hierarchical data structures”**



# What does it mean for samples to be independent of each other?



**How can we know for SURE the difference between A, B, and C is due to treatment and not due to difference between flasks?**

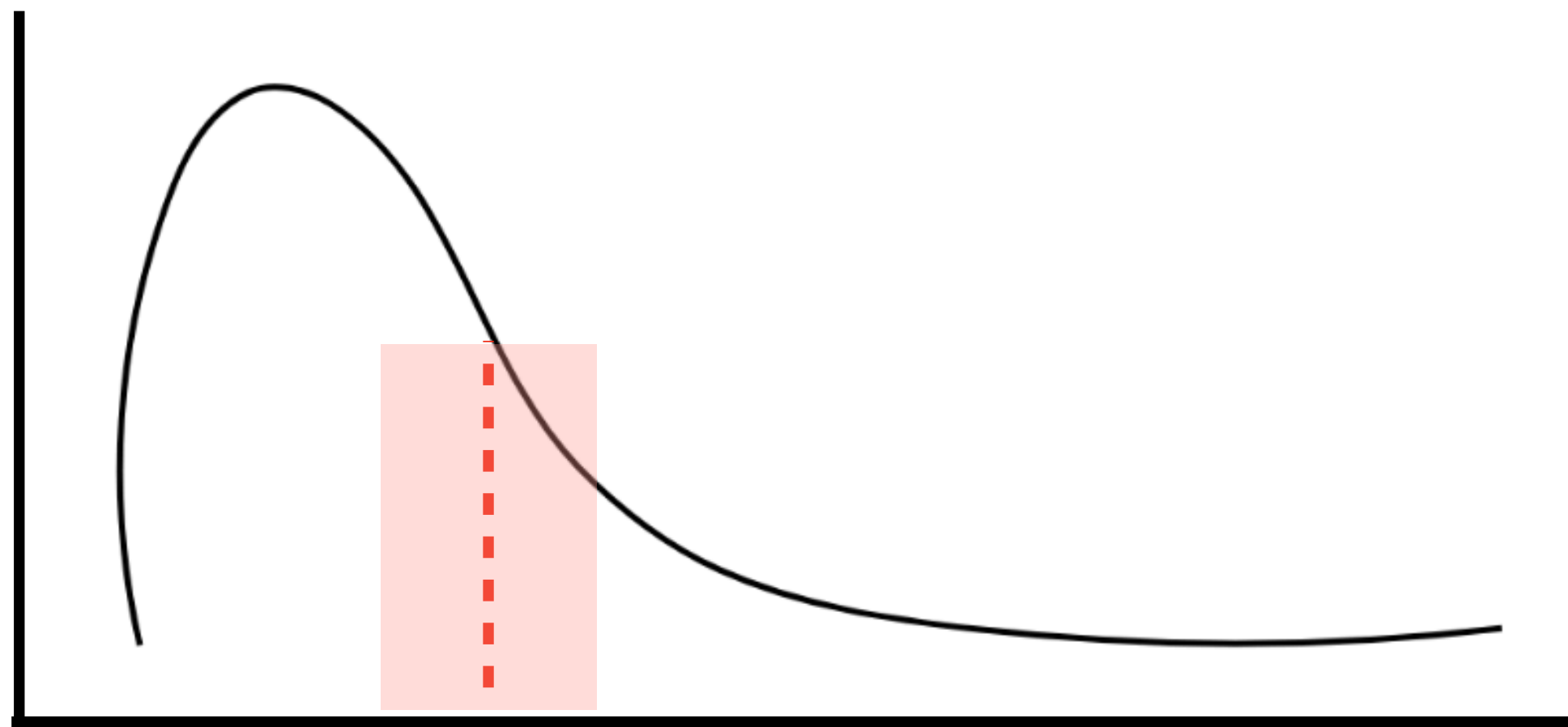
**“Hierarchical data structures”**

# Assumptions for estimating confidence intervals

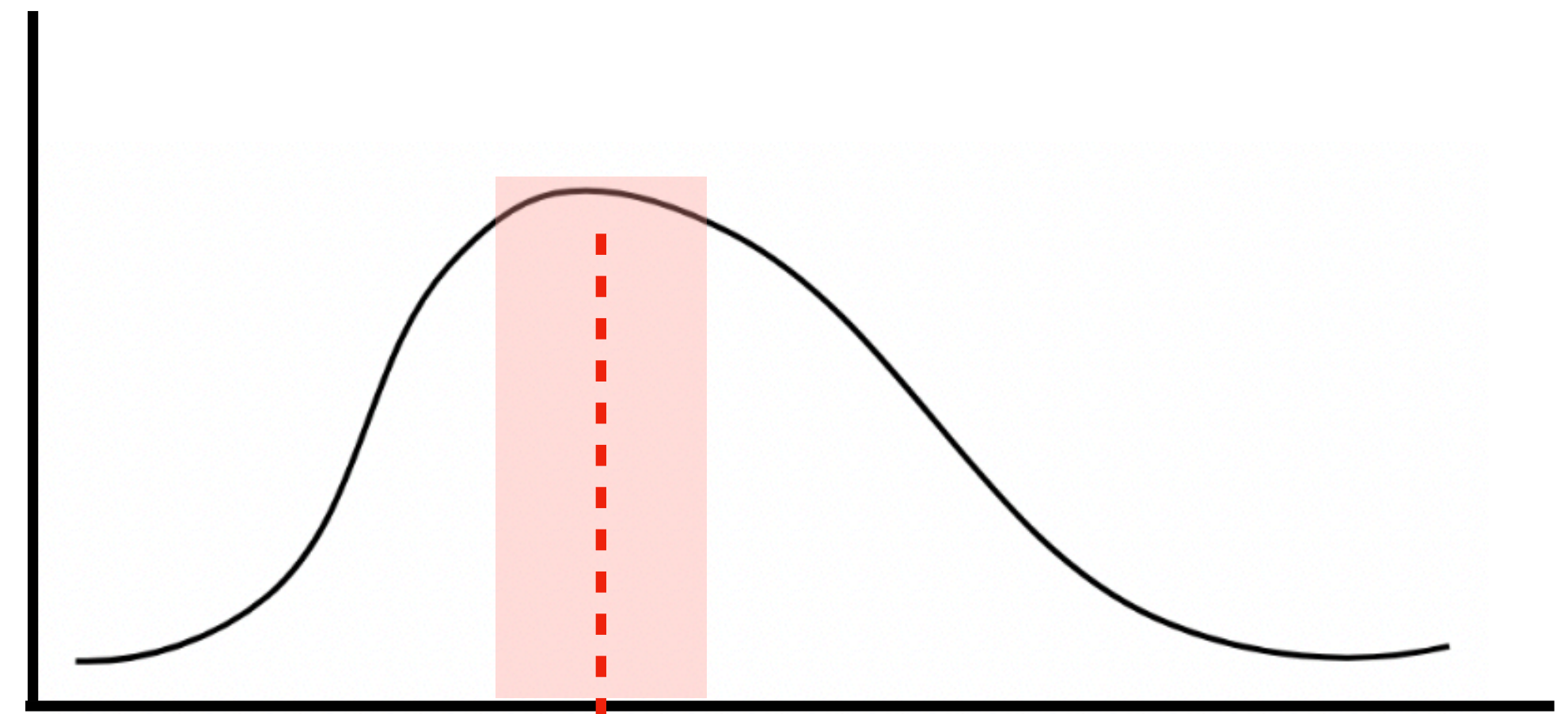
- **Conditions on the design of the study:**
  - (1) Data is a random sample from a large population
  - (2) Observations in the sample must be independent of each other
- **Conditions on the form of the population distribution**
  - (3) If  $n$  is small, the population distribution must be ~normal
  - (4) If  $n$  is large, the population distribution doesn't have to be normal



# Why does it matter if our population is normally distributed?



**Is the mean even a meaningful measure for this population?**



**If  $n$  is large enough, the sample distribution will be normal regardless of the shape of the population**

# How can we tell if a population is normally distributed?

- **Plot distribution (i.e. histogram)**
  - *Every analysis should begin with an inspection of the data and the points that lie far from the center*
- **Quantile plot**
- **Shapiro-wilks test for non-normality**
- **If not normal distribution, try a data transformation**

# Announcements

- Extra practice problems from the textbook posted to GitHub and Canvas
  - Note: no solutions, but use your classmates/TAs for help!
- If you need help keeping track of the stats R functions we have been using, check out the `stats_R_cheatsheet.md` on GitHub!