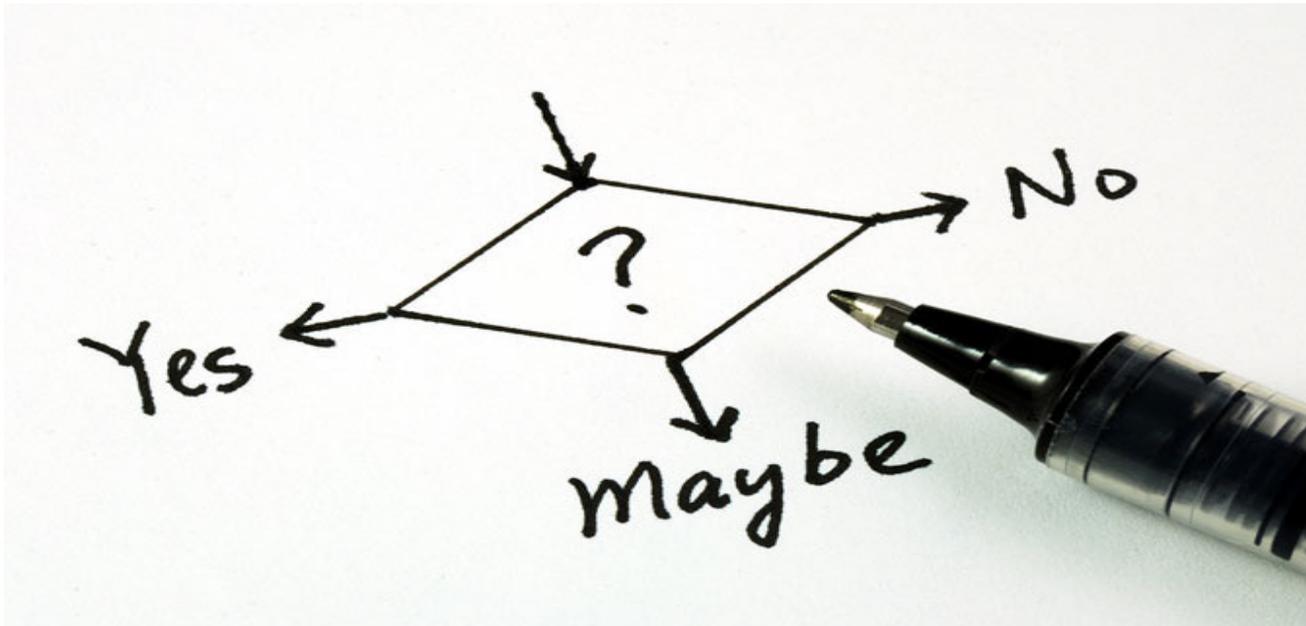


Don't hesitate to raise your hand during lectures if you have any questions.

I also try to “read the room” and will step in when I sense that students may be unsure or have questions they haven't yet voiced.



Statistics in the science of supporting decision-making under incomplete information / knowledge



We can't measure everything:
statistics are based on samples!

Biologists are relatively small: we collect smaller
number of things to generalize to all things!



- From Chris Lortie

Mapping tree density at a global scale

T. W. Crowther¹, H. B. Glick¹, K. R. Covey¹, C. Bettigole¹, D. S. Maynard¹, S. M. Thomas², J. R. Smith¹, G. Hintler¹, M. C. Duguid¹, G. Amatulli³, M.-N. Tuanmu³, W. Jetz^{1,3,4}, C. Salas⁵, C. Stam⁶, D. Piotta⁷, R. Tavani⁸, S. Green^{9,10}, G. Bruce⁹, S. J. Williams¹¹, S. K. Wiser¹², M. O. Huber¹³, G. M. Hengeveld¹⁴, G.-J. Nabuurs¹⁴, E. Tikhonova¹⁵, P. Borchardt¹⁶, C.-F. Li¹⁷, L. W. Powrie¹⁸, M. Fischer^{19,20}, A. Hemp²¹, J. Homeier²², P. Cho²³, A. C. Vibrans²⁴, P. M. Umunay¹, S. L. Piao²⁵, C. W. Rowe¹, M. S. Ashton¹, P. R. Crane¹ & M. A. Bradford¹

A study led by Yale University researchers has found that there are over 3 trillion trees on Earth - but they are **disappearing at an alarming rate**.

The study found that there are around 3.04 trillion trees on Earth, or around 422 for each person on the planet.

The number is a huge increase on the previous global estimate, which was just over 400 billion trees worldwide.

The study was based on on-the-ground data about the number of trees in more than 400,000 plots of forest from all continents except Antarctica.

Source - <https://www.independent.co.uk/environment/how-many-trees-are-there-on-earth-10483553.html>



Statistics are based on samples!

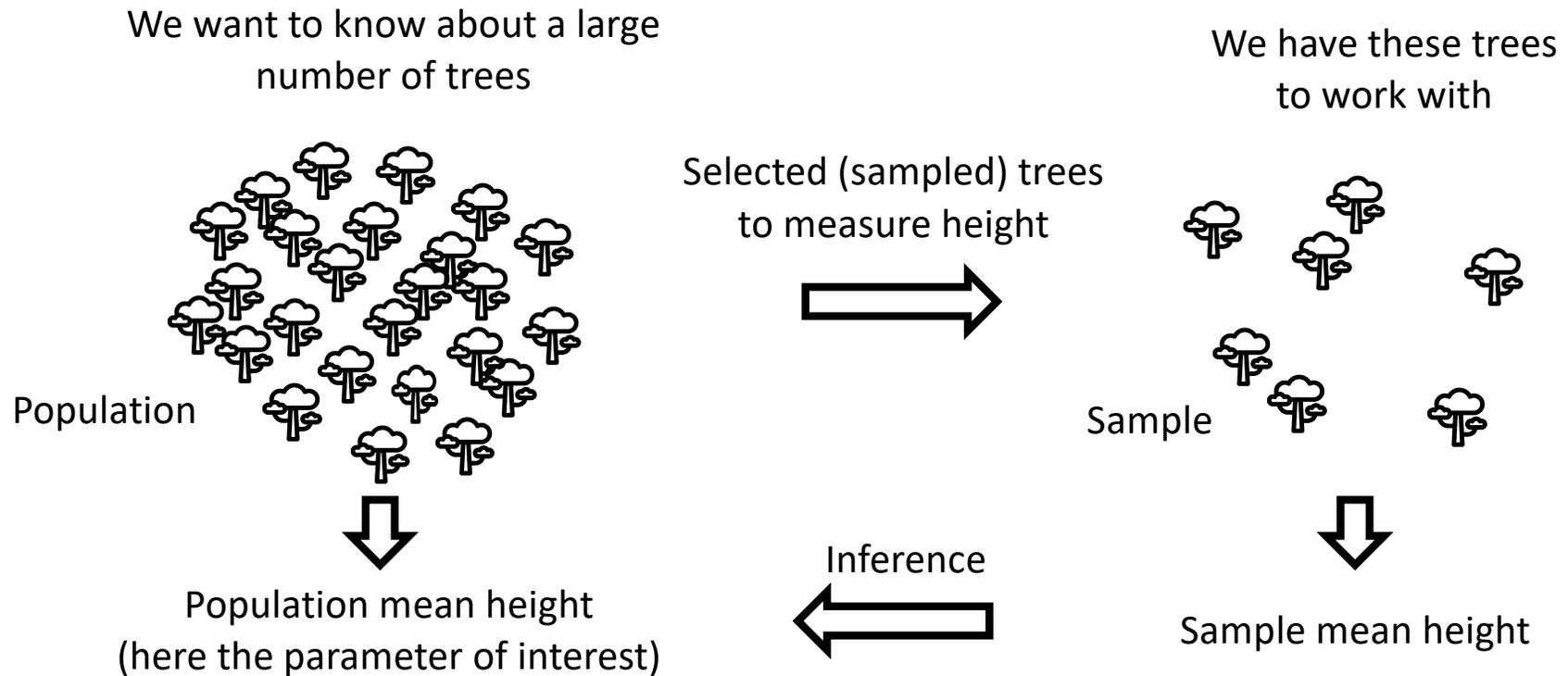
A central goal of statistics is to learn about an entire population (e.g., the average height of a plant species) using information from a sample - a smaller subset of observations or individuals drawn from that population.

e.g., = *exemplī grātiā* (for the sake of example!)

— Biologists are relatively small



A central goal of statistics is inference: drawing conclusions about unknown population quantities (e.g., average height of a tree species in a forest) from **sample** data.



A sample is a subset of observations or individuals drawn from a larger population, used to learn about that population.

Some jargon is key to optimize communication & understand concepts more clearly

What is a *population*, a *sample*, an *observation*, a *variable*, a *parameter*?

Jargon gets an unfair treatment but learning and working in most fields would be very difficult without it.

Jargon is a real time saver!



Statistical inference is based on samples!

The central goal of statistics is inference: drawing conclusions about unknown population quantities (e.g., average height of a tree species in a forest) from sample data.

Infer an unknown quantity =
use a sample to produce information (e.g., an average value) about a chosen statistical population of interest (e.g., trees in a forest; potato-chip bags produced by a factory in a year).

Statistical (important) definitions: **POPULATION**

A population is the entire collection of individual units (or observation units) that share a common property or set of properties.

It is from this group that we aim to generalize knowledge about unknown quantities (e.g., mean of the population) based on a subset of individual units, known as a sample.

Examples -

- Coffee drinkers in Quebec.
- Coffee drinkers in Canada.
- Coffee drinkers in Canada that run in the morning.

Statistical (important) definitions: **POPULATION**

A population is the entire collection of **individual units** (or **observation units**) that share a common property or set of properties. It is from this group that we aim to generalize knowledge about unknown quantities based on a subset of individual units (observation), known as a **sample**.

What is the average height and average weight of coffee drinkers in Canada that run in the morning?

Individual unit (or observation unit) = someone living in Canada that drinks coffee and runs in the morning.

Statistical (important) definitions: **POPULATION**

A population is the entire collection of **individual units** (or **observation units**) that share a common property or set of properties. It is from this group that we aim to generalize knowledge about unknown quantities based on a subset of individual units (observation), known as a **sample**.

What is the average height and average weight of coffee drinkers in Canada that run in the morning?

Individual unit (or observation unit) = someone living in Canada that drinks coffee and runs in the morning.

Properties = Live in Canada, drink coffee and run in the morning.

Statistical (important) definitions: **POPULATION**

A population is the entire collection of **individual units** (or **observation units**) that share a common property or set of properties. It is from this group that we aim to generalize knowledge about unknown quantities based on a subset of individual units (observation), known as a **sample**.

What is the average height and average weight of coffee drinkers in Canada that run in the morning?

Individual unit (or observation unit) = someone living in Canada that drinks coffee and runs in the morning.

Properties = Live in Canada, drink coffee and run in the morning.

Observation (or data point) = set of one or more quantities (measurements) on a single observation unit; ex. the weight and height of someone living in Canada that drinks coffee and run in the morning.

Statistical (important) definitions: **POPULATION**

A population is the entire collection of **individual units** (or **observation units**) that share a common property or set of properties. It is from this group that we aim to generalize knowledge about unknown quantities based on a subset of individual units (observation), known as a **sample**.

What is the average height and average weight of coffee drinkers in Canada that run in the morning?

Individual unit (or observation unit) = someone living in Canada that drinks coffee and runs in the morning.

Properties = Live in Canada, drink coffee and run in the morning.

Observation (or data point) = set of one or more quantities (measurements) on a single observation unit; ex. the weight and height of someone living in Canada that drinks coffee and run in the morning.

Sample = subset of observation units from all possible observations in the population.

A **sample** of 11 individuals from the target **population** (**PROPERTIES**: Canadians that drink coffee and run in the morning); target population = population of interest.

Individual unit (or observation unit) = a Canadian that drinks coffee and runs in the morning.



Individual	Weight (kg)	Height (cm)
1	75.5	172
2	55.3	152
3	61.2	164
4	50.3	148
5	99.4	192
6	66.2	171
7	75.3	169
8	74.6	182
9	60.5	162
10	93.5	184
11	73.6	169

Observation (or data point) = set of one or more quantities (measurements) on a single observation unit; ex. the weight and height of someone living in Canada that drinks coffee and run in the morning.

TWO different observations are in red squares.

Statistical (important) definitions: **POPULATION**

The size of a population is often unknown. For example, we may not know how many people in Canada drink coffee and go for a run in the morning.

In many cases, the population can be so large that it is considered 'infinite' for practical purposes (we will explore this concept further later in the semester).

Additionally, populations are dynamic, meaning they can change over time.

Mapping tree density at a global scale

T. W. Crowther¹, H. B. Glick¹, K. R. Covey¹, C. Bettigole¹, D. S. Maynard¹, S. M. Thomas², J. R. Smith¹, G. Hintler¹, M. C. Duguid¹, G. Amatulli³, M.-N. Tuanmu³, W. Jetz^{1,3,4}, C. Salas⁵, C. Stam⁶, D. Piotta⁷, R. Tavani⁸, S. Green^{9,10}, G. Bruce⁹, S. J. Williams¹¹, S. K. Wiser¹², M. O. Huber¹³, G. M. Hengeveld¹⁴, G.-J. Nabuurs¹⁴, E. Tikhonova¹⁵, P. Borchardt¹⁶, C.-F. Li¹⁷, L. W. Powrie¹⁸, M. Fischer^{19,20}, A. Hemp²¹, J. Homeier²², P. Cho²³, A. C. Vibrans²⁴, P. M. Umunay¹, S. L. Piao²⁵, C. W. Rowe¹, M. S. Ashton¹, P. R. Crane¹ & M. A. Bradford¹

A study led by Yale University researchers has found that there are over 3 trillion trees on Earth - but they are **disappearing at an alarming rate**.

The study found that there are around 3.04 trillion trees on Earth, or around 422 for each person on the planet.

The number is a huge increase on the previous global estimate, which was just over 400 billion trees worldwide.

The study was based on on-the-ground data about the number of trees in more than 400,000 plots of forest from all continents except Antarctica.

Source - <https://www.independent.co.uk/environment/how-many-trees-are-there-on-earth-10483553.html>



Let's pause and think:

If the number of trees is revised upward (from 400 billion to 3 trillion), should we automatically conclude that earlier estimates of the average tree size were wrong?



Let's pause and think:

If the number of trees is revised upward (from 400 billion to 3 trillion), should we automatically conclude that earlier estimates of the average tree size were wrong?

Not necessarily. A higher population size changes **how many trees exist**, but it doesn't automatically change **the average size** - that depends on *which trees* are being counted and how the average was estimated.



Statistical populations *versus* biological populations (let's not mix the two)

In biology, a population refers to all organisms of the same group or species that live in a specific geographical area.

In statistics, a population is a set of similar items, whether living or non-living, that are of interest for answering a particular research question.

For clarity, it's best to refer to it as a 'statistical population' when the goal is to infer quantities from it.

Let's take a break – 1 minute



Statistics is the science of supporting decision-making under incomplete information, most often by analyzing samples drawn from populations of unknown size.

In other words: Statistics is the science of inference and decision-making under uncertainty, using samples to learn about populations we cannot fully observe.



Statistical (important) definitions: **POPULATION**

Examples -

Stars in the sky (“infinite”)

Sand in a river (“infinite”)

Countries in Europe (finite)

Bags of potato chips in a factory (finite)

Again, to avoid confusion, we often use the term “statistical population” instead of just “population”; and some refer to “statistical universe”.

What we want to know determines what we measure and which statistical population we are studying.

Question	Observational unit	Statistical population
What proportion of the plants are flowering?	An individual plant	All the plants in the ecological population
How many seeds per flower?	An individual plant in flower	All the plants in flower
How many seeds per white-flowered plant?	An individual white-flowered plant	All the white-flowered plants in flower
How many plants/m ² in the field?	An area of m ²	All the areas of m ² in the field
How long are the stamens?	A stamen	All the stamens
How much time do bees spend on a visit to a flower?	A visit by a bee to a flower	All the visits made by bees to flowers
How many bees visit in a 5-minute observation period?	A 5-minute observation period	All the 5-minute observation periods which could be made

A **parameter** is a quantity describing a statistical population, whereas an **estimate** (or statistic) is the same quantity but calculated from a sample.

What is the average height of trees
across all species worldwide?

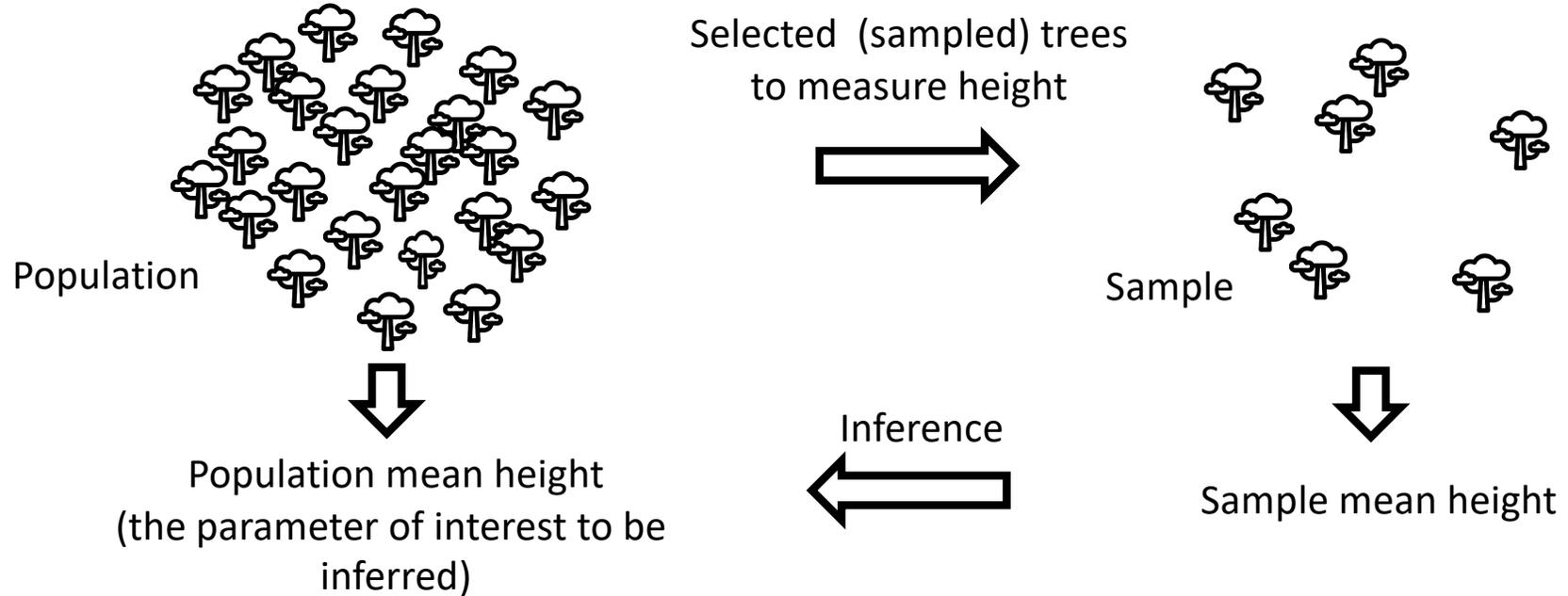
The parameter of interest (unknown) is the average height of all trees.

The estimate of interest (known) is the average height of a smaller group of trees (sample).

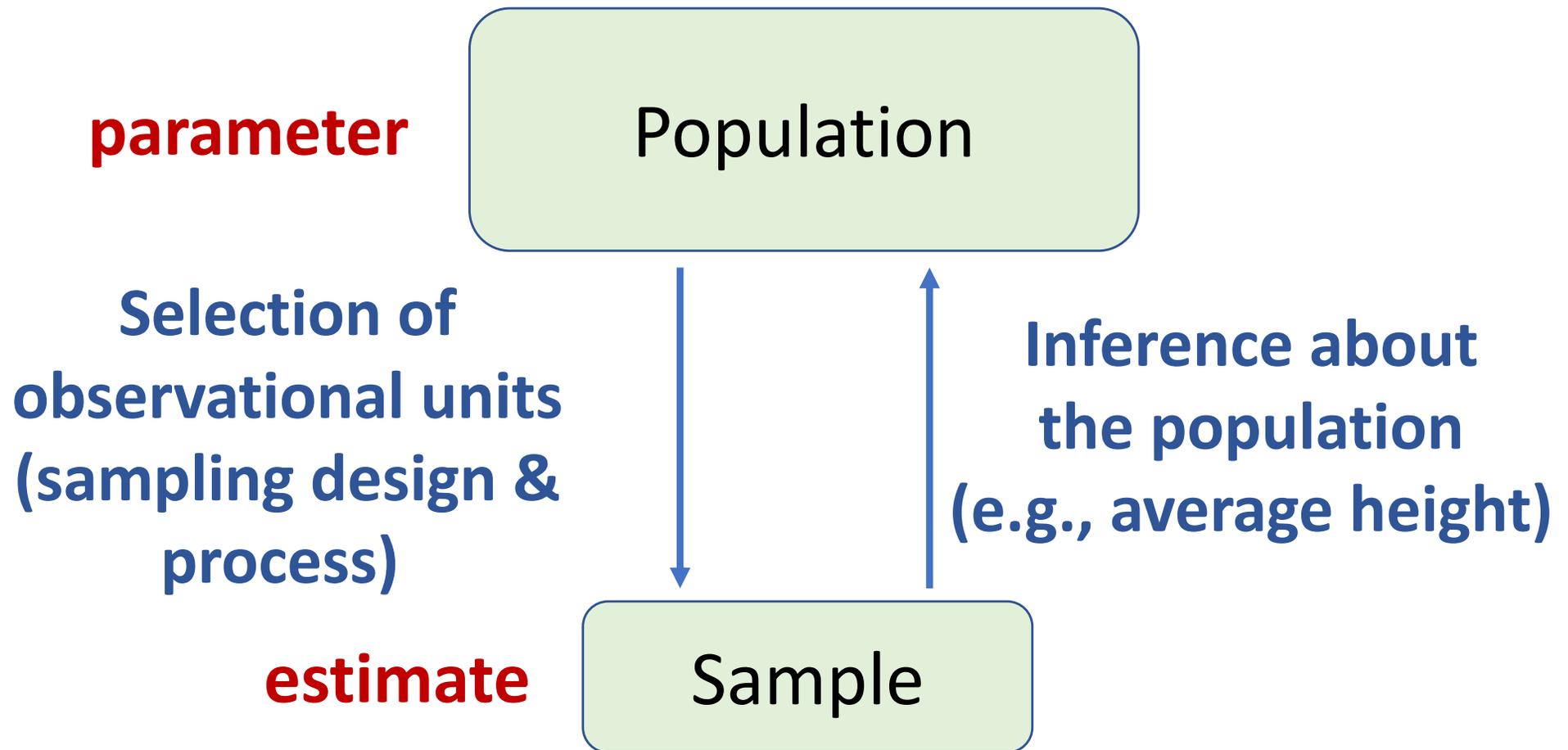
A central goal of statistics is inference: drawing conclusions about unknown population quantities (e.g., average height of a tree species in a forest) from **sample** data.

We want to learn about these trees (which, in most cases, we don't even know how many there are).

We have these trees to work with



A central goal of statistics is inference: drawing conclusions about unknown population quantities (e.g., average height of a tree species in a forest) from **sample** data.



Variables!

A *variable* is any characteristic, number, or quantity that can be measured or counted and varies among observation units. Examples of variables include height, weight, age, gender, and eye color.

Recognizing the type of *variable* is crucial, as it often determines the appropriate type of statistical analysis.

Variables (e.g., height, biomass) differ among observation units (e.g., individual trees).



A **sample** of 11 individuals from the target populations (**PROPERTIES**: Canadians that drink coffee and run in the morning).

Individual unit (or observation unit) = a Canadian that drinks coffee and runs in the morning.



VARIABLES



Individual	Weight (kg)	Height (cm)
1	75.5	172
2	55.3	152
3	61.2	164
4	50.3	148
5	99.4	192
6	66.2	171
7	75.3	169
8	74.6	182
9	60.5	162
10	93.5	184
11	73.6	169



An **observation unit** (or simply observation) contains the values of all **variables** of interest, such as the height and weight of a Canadian who drinks coffee and runs in the morning.

TWO different observations.



Types of variables

CATEGORICAL VARIABLES - describe membership in a category or group; characteristics of observations that do not have magnitude on a numerical scale. They can be:

Nominal (name):

- Survival (alive or dead),
- Method of disease transmission (e.g., water, air, animal vector),
- Eye colors (amber, blue, brown, gray, green, hazel, or red),
- Breed of a dog (e.g., collie, shepherd, terrier).

or Ordinal (ordered):

- Life stage (e.g., egg, larva, juvenile, adult),
- Snake bite severity score (e.g., minimal, moderate, severe),
- Size class (e.g., small, medium, large).

The case of MONTHS and WEEKDAYS as categorical variables [revised slide]

These variables (e.g., Day of the week or month of bird observation) describe time-related categories, and how they are treated statistically depends on the question being asked.

Weekdays and months have a natural (cyclical) order (Monday → Sunday; January → December), so they are technically ordinal variables.

However, the numerical distance between days is not inherently meaningful; coding Monday = 1 and Tuesday = 2 does not create a meaningful “distance” between them, i.e., the difference between Monday and Tuesday is not intrinsically the same kind of quantity as the difference between Tuesday and Friday.

If one wants to quantify time differences, a continuous measure such as hours (or days) should be used instead.

For this reason, weekdays and months are often treated as *categorical variables* in practice - either *nominal* (when only category matters) or *ordinal* (when their ordering is relevant to the analysis).

Types of variables [**revised slide**]

NUMERICAL VARIABLES - characteristics of observations have magnitude on a numerical scale.

Continuous (can take any real-number value)

- Core body temperature (e.g., degrees Celcius, °C),
- Territory size of a bird (e.g., hectares),
- Size of fish (e.g., cm)

Discrete (only take indivisible units though they can be non-integers, i.e., cannot be subdivided into smaller, meaningful units while maintaining their discrete nature)

- Age at death (e.g., years),
- Number of amino acids in a protein,
- Number of eggs in a bird nest.
- *Confusing but true*: 0.5 mg/mL, 1.0 mg/mL, and 1.5 mg/mL. These data are discrete because the researcher only used these specific, predetermined concentrations, not any value in between.

The case of YEARS as a numeric variable [new slide]

Years (e.g., year of bird observation) describe time-related values, and how they are treated statistically depends on the question being asked.

Years follow a natural, linear order (... → 2022 → 2023 → 2024 → 2025), and - unlike months or weekdays - the spacing between years is consistent.

The numerical distance between years is inherently meaningful: the difference between 2020 and 2021 represents the same amount of time as the difference between 2023 and 2024.

For this reason, years are typically treated as a numeric (interval) variable, and differences between years can be used directly to quantify time change (e.g., trends, rates, or durations).

NOTE: In contrast, “day of the month” cannot generally be as a numeric variable because months have different lengths and the scale is cyclical; it is therefore usually treated as a categorical (ordinal or nominal) variable unless converted into a continuous measure such as day of year.

NUMERICAL VARIABLES [new slide]

Length, Numbers Between an interval & Cardinality

Length (or distance) - a measurement of change = $|b - a|$

Examples: Length from 1 to 2 $\rightarrow 2 - 1 = 1$ Length from 2 to 2 $\rightarrow 2 - 2 = 0$

Length tells you *how far apart two numbers are on the number line*, not how many numbers are between them.

Length is the most important in statistics; it measures how much something changes, and length (distance) is what measures change, whereas “how many numbers are between” and “cardinality” mostly describe how we label or count values.

Discrete variables (e.g., counts, number of trees, number of students) - length corresponds to how many units have changed. Example:
2 trees \rightarrow 5 trees; length = 3 \rightarrow an increase of 3 individual trees

Continuous variables (e.g., time, temperature, height). Length corresponds to how much change occurred on a continuous scale. Example:
2 seconds \rightarrow 5 seconds; length = 3 \rightarrow 3 seconds of time passed

There are infinitely many possible intermediate values for continuous variables, but that does not change the meaning of length.

NUMERICAL VARIABLES [new slide]

Length, Numbers Between an interval & Cardinality

Numbers between an interval: it depends on whether numbers are discrete or continuous. “numbers between” is usually irrelevant statistically (important but no need to know for BIOL-322).

Discrete variables (i.e., discrete numbers) and for an interval $[a,b]$, the formula is:
Number between = $\max(0, b - a - 1)$

Examples:

$$[1,2] \rightarrow 2 - 1 - 1 = 0$$

$$[2,2] \rightarrow 2 - 2 - 1 = -1 \rightarrow 0; \text{ remember } \max(0, b - a - 1)$$

$$[2,4] \rightarrow 4 - 2 - 1 = 1 \rightarrow (\text{only number } 3)$$

Note again: a discrete variable doesn't have to be made of integer numbers.

Continuous variables (i.e., real numbers)

Examples:

$$[1,2] \rightarrow \text{infinitely many real numbers}$$

$$[2,3] \rightarrow \text{infinitely many real numbers}$$

$$[2,2] \rightarrow \text{zero, because the length is } 0$$

With continuous numbers, “having space” (positive length) automatically means infinitely many numbers inside.

NUMERICAL VARIABLES [new slide]

Length, Numbers Between an interval & Cardinality

Cardinality of an interval: the number of elements in a set and depends on whether numbers are discrete or continuous. “cardinality” is usually irrelevant statistically (important but no need to know for BIOL-322).

Discrete variables (i.e., discrete numbers) and for an interval $[a,b]$, the formula is: Number between = $\max(0, b - a - 1)$

Examples:

Set $\{1,2\}$ → cardinality = 2; discrete intervals are contained by $\{ \}$ - curly brackets

Set $\{2,3,4\}$ → cardinality = 3

Continuous variables (i.e., real numbers)

Examples:

$(1,2)$ → infinitely many real numbers → infinite cardinality

$(2,3)$ → infinitely many real numbers → infinite cardinality

If the interval has positive length, like $(1,2)$ or $(2,3)$, its cardinality is infinite.

Confusing but true: The only case where cardinality is not infinite in the continuous world is when the interval has zero length, like:

$[2,2]$ → length 0 → cardinality = 1: interval closed ($[]$) at both ends, but the ends coincide.

$(2,2)$ → length 0 → cardinality = 0 (empty set): interval open ($()$) at both ends, , and because the endpoints coincide, nothing is included.

NUMERICAL VARIABLES [new slide]

THE CASE OF TIME: Length, Numbers Between an interval & Cardinality

Length (what matters most in statistics)

Length = $|b-a|$: measures **change in time**, expressed in a **unit** (e.g., seconds).
 $2\text{ s} \rightarrow 5\text{ s} \Rightarrow \text{length} = \mathbf{3\text{ seconds}}$ (real duration).

(*Physics note*): below **Planck time** ($\sim 5.4 \times 10^{-44}\text{ s}$), time may not be divisible; “continuous time” is a mathematical idealization; in physics, time is not assumed to be divisible forever.

Numbers between (depends on discrete vs continuous)

Discrete seconds: number between = $\max(0, b-a-1)$.
Between 2 and 5 s \rightarrow 2 (3 and 4).

Continuous seconds: any interval with positive length contains infinitely many real numbers in seconds.

Between 2 and 5 s \rightarrow infinitely many values (2.1, 2.01, ...).

What is infinite is the **number of numerical labels in seconds**, **not** the time (still 3 s).

Cardinality (size of the set)

Discrete interval in seconds: $\{2,3,4,5\} \rightarrow$ cardinality = 4.

Continuous interval in seconds: $(2,5) \rightarrow$ infinite cardinality; $[2,2] \rightarrow 1$; $(2,2) \rightarrow 0$.

Discrete versus Continuous Probability [new slide]

Discrete variables: it's possible to assign a probability to a specific value in a series of discrete numbers.

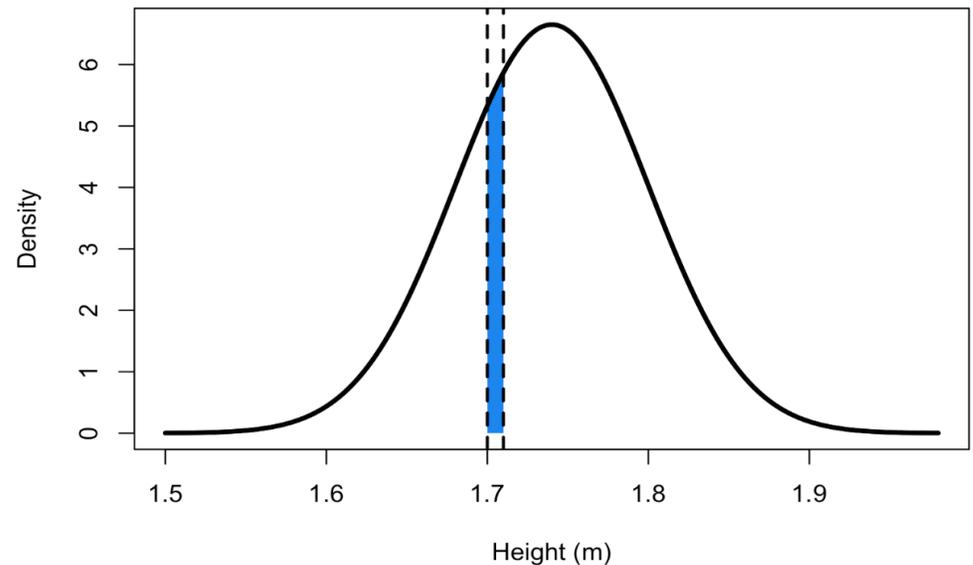
Continuous variables: the probability of any exact value is zero:

The probability that someone is exactly 1.700000000000... m is 0..).

Asking for the probability that someone is exactly 1.700000... m tall is like asking for the probability of picking *one specific grain of sand* from all the sand on a beach - it's so small that we treat it as zero.

Instead, probabilities are defined over intervals between two values. In statistics, this means that for continuous variables we ask questions like: *what is the probability that someone's height lies between 1.70 m and 1.71 m?*

Since there are infinitely many possible heights between 1.70 m and 1.71 m, probability gets spread so thinly across them that the chance of any one exact height is effectively zero — and in any case too small to be useful.



Statistical variables – important note

Variables are defined by their type (e.g., height, length), not by their measurement units (e.g., cm).

As such, arm length and leg length can be both measured in centimeters, but they are TWO different variables.

Enjoy your definitions!

