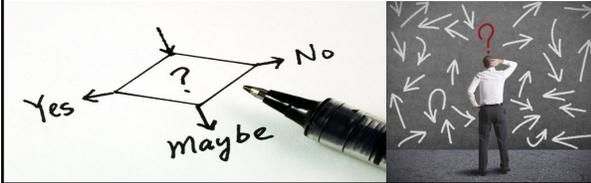


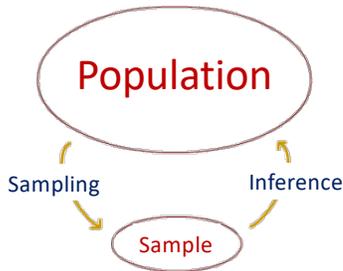
Lecture 7: estimating & making inferences with uncertainty – samples and biases

The science of aiding decision-making with incomplete information (i.e., without certainty)



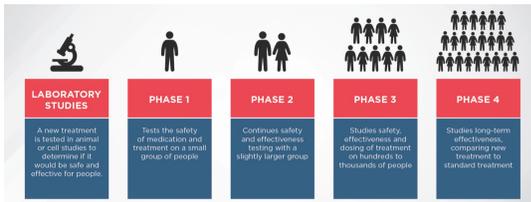
1

Inferential process



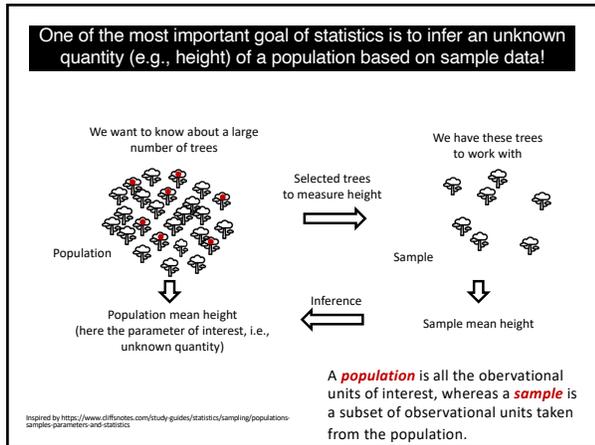
2

A good example of sampling: Stages of Clinical Trials

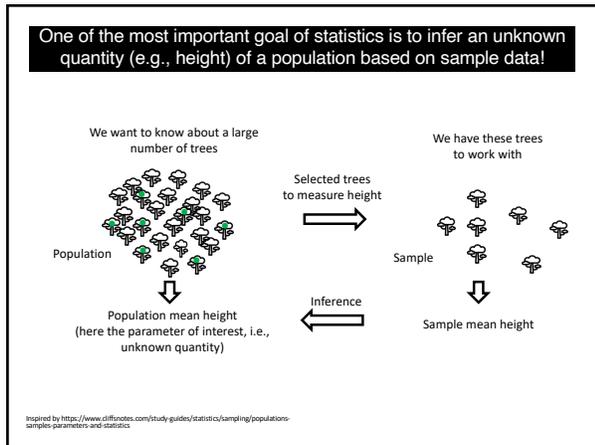


<https://news.sanfordhealth.org/coronavirus-disease-2019-covid-19/coronavirus-expert-q-a/sanford-research-vaccine-phases/>

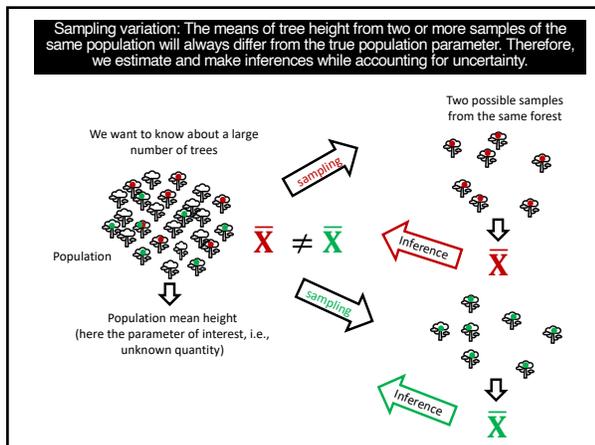
3



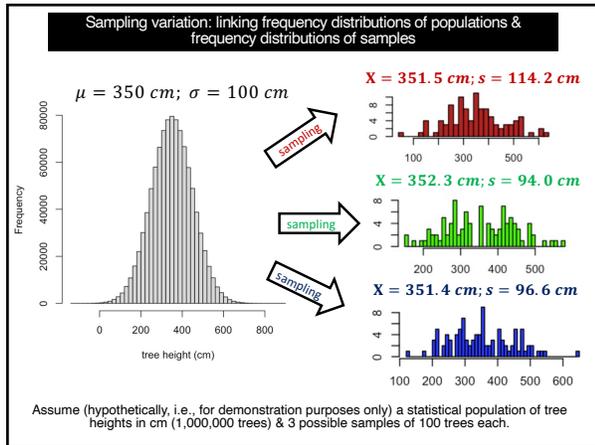
4



5



6



7

Sampling variation: linking frequency distributions of populations & frequency distributions of samples

Assume (hypothetically, i.e., for demonstration purposes only) a statistical population of tree heights in cm (1000000 trees) & 3 possible samples of 100 trees each.

How many possible samples of 100 trees out of a population with 1000000 trees?

10768272362e+432 (zeros)

For comparison: the **human body** consists of about 37.2 trillion **cells** (3.72e+13 zeros)

8

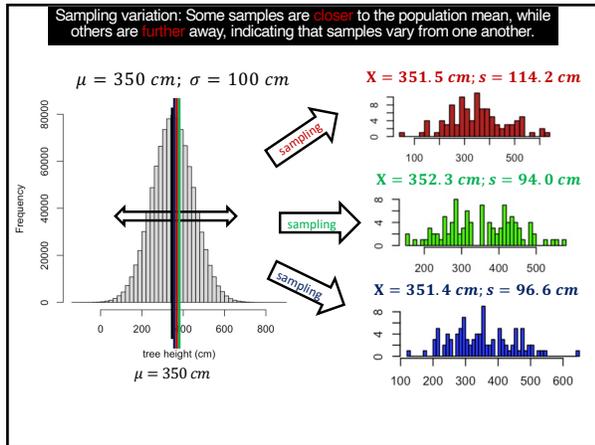
In most real-world studies, however, we typically take only one sample from the target statistical population.

1,000,000 trees 100 trees

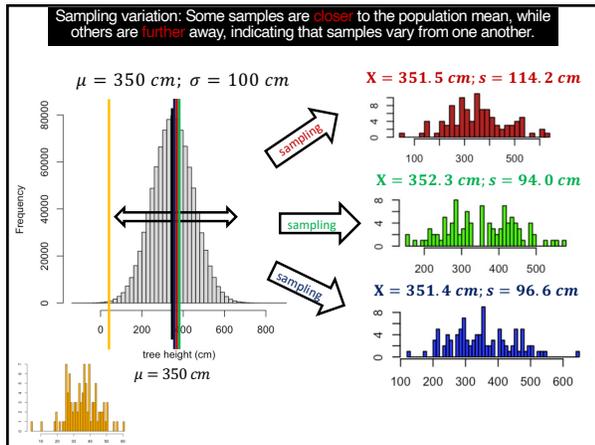
Critical understanding:

- [1] Values of descriptive statistics from samples are never exactly the same as the corresponding population values due to sampling error.
- [2] However, this does not mean that inferences drawn from samples are incorrect (we'll explore this in more detail later). Sample values serve as good approximations of the true population values.
- [3] These approximations can range from good (when the sample statistic is close to the true population value) to poor (when it is far from the true value).

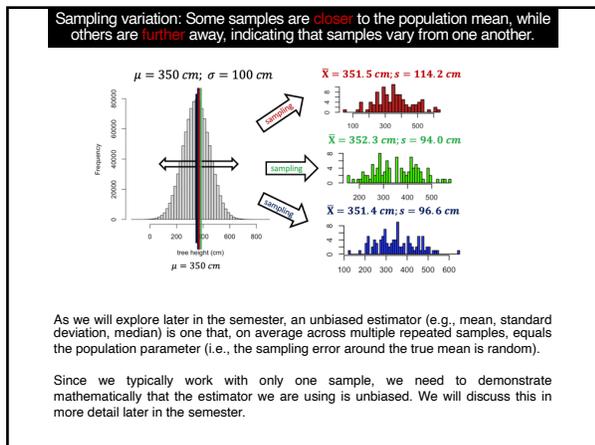
9



10



11



12

In most real-world studies, however, we typically take only one sample from the target statistical population

$\mu = 350 \text{ cm}; \sigma = 100 \text{ cm}$

$\mu = 350 \text{ cm}$

$X = 352.3 \text{ cm}; s = 94.0 \text{ cm}$

sampling →

Critical understanding:

[3] Again, approximations can be either good (when the sample value is close to the true population value) or bad (when the sample value is far from the true value). Later, you will understand why we use terms like "close" and "far" to describe how samples relate to their populations.

[4] To build confidence, it would be helpful to estimate the expected margin of error.

13

How wrong one could be in trusting their sample values to estimate the population value (i.e., parameter)?

$\mu = 350 \text{ cm}; \sigma = 100 \text{ cm}$

$\mu = 350 \text{ cm}$

$X = 352.3 \text{ cm}; s = 94.0 \text{ cm}$

sampling →

←

Critical understanding:

[5] As we will see later, the variation among observations within a single sample (standard deviation) can help estimate how far sample means might be from the true population mean, providing an idea of potential error.

[3] Again, approximations can be either good (when the sample value is close to the true population value) or bad (when the sample value is far from the true value). Later, you will understand why we use terms like "close" and "far" to describe how samples relate to their populations.

[4] To build confidence, it would be helpful to estimate the expected margin of error.

14

Key concepts underlying statistics and statistical thinking

- Uncertainty: never being able to determine the true population parameter.
- The risk of error arises when decisions are based on estimates. If the estimate is close to the true value, the impact is minimal. However, when the estimate is far from the true value, decisions may lead to significant errors.
- Assessing uncertainty (error/risk) then becomes crucial.
- Sample variability - Answer may change with different sample data.
- Accuracy (how close to the true population value).

15

**Key concepts:
Statistics is based on samples!**

Sample quantities (e.g., mean, median, standard deviation, interquartile range) almost always vary from sample to sample, introducing some level of uncertainty.

Therefore, any estimates or inferences we make about population values inherently involve some degree of uncertainty, both in terms of variation between samples and the difference between sample estimates and the true population value.

As we will explore later, the variation among observations within a sample (standard deviation) can help estimate how far sample means may deviate from the true population mean, allowing us to assess uncertainty and associated risks.

16

Don't forget to watch all the material in our WebBook.

Understanding sampling variation with dance.



<https://www.youtube.com/watch?v=5fGu8hvdZ6s>

17

Let's take a break - 1 minute



18

Random sampling helps minimize sampling error (i.e., how close or far sample values are from the true population value for the statistic of interest) and reduces inferential bias.

The key requirement for the methods presented in this course (and in statistics in general) is that the data come from a random sample. A random sample must meet two essential criteria:

- 1) Every observational unit in the population (e.g., individual tree) have an **equal chance** of being included in the sample.
- 2) The selection of observational units in the population (e.g., individual tree) must be **independent**, i.e., the selection of any unit (e.g., individual tree) of the population must not influence the selection of any other unit.

➔ Samples are biased when certain observational units in the target population have a lower or higher probability of being selected.

19

Before I forget!!!!

2) The selection of observational units in the population (e.g., individual tree) must be **independent, i.e.**, the selection of any unit (e.g., individual tree) of the population must not influence the selection of any other unit.

i.e. = id est (it is)

20

Properties of samples - Precision and accuracy:
the major goal of sampling is to increase accuracy & precision

Precise

Tree population
 $\mu = 350 \text{ cm}$

1 000 000 trees (heights in cm)

Accurate

4 random samples of 100 trees each

The average sample mean shows little variation from sample to sample (precision) and, on average, is similar to the population mean (accuracy).

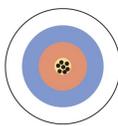
These 4 samples are precise and accurate

21

Properties of samples - Precision and accuracy:
the major goal of sampling is to increase accuracy & precision

Precise

Accurate



Imagine the bull's eye as the population parameter (in this case, the mean tree height), and the points as possible sample mean values for tree height (i.e., estimates).

Accurate = sample values (e.g., sample means) tend to be close to the true population value.

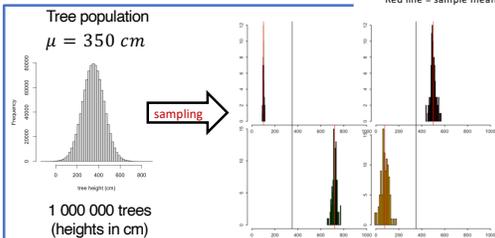
Precise = Sample values (e.g., sample means) tend to be similar to each other, regardless of whether they are close to or far from the true population value.

22

Properties of samples - Precision and accuracy:
the major goal of sampling is to increase accuracy & precision

Accurate

Imprecise



Tree population
 $\mu = 350\text{ cm}$

1 000 000 trees
(heights in cm)

4 random samples
of 100 trees each

The average sample mean varies considerably from sample to sample (imprecise), but on average, it is close to the population mean (accurate). In other words, the average of these four samples is very close to the true population value.

Black line = Population mean
Red line = sample mean

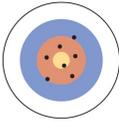
These 4 samples are imprecise but accurate

23

Properties of samples - Precision and accuracy:
the major goal of sampling is to increase accuracy & precision

Imprecise

Accurate

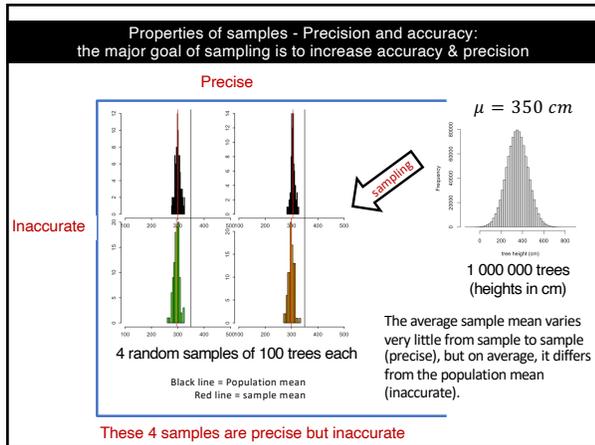


Imagine the bull's eye as the population parameter (in this case, the mean tree height), and the points as possible sample mean values for tree height (i.e., estimates).

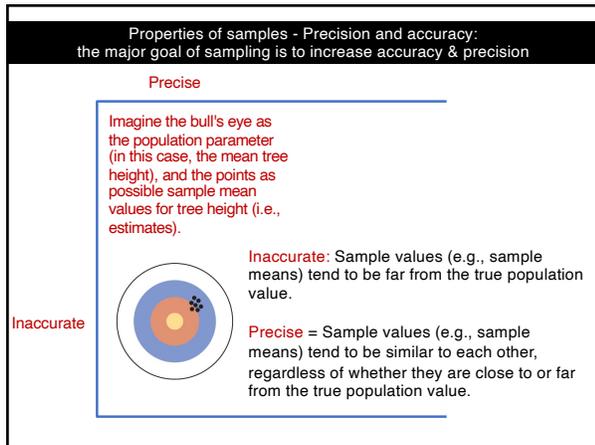
Accurate = sample values (e.g., sample means) tend to be close to the true population value.

Imprecise: Sample values (e.g., sample means) tend to vary widely from each other, regardless of whether they are close to or far from the population value.

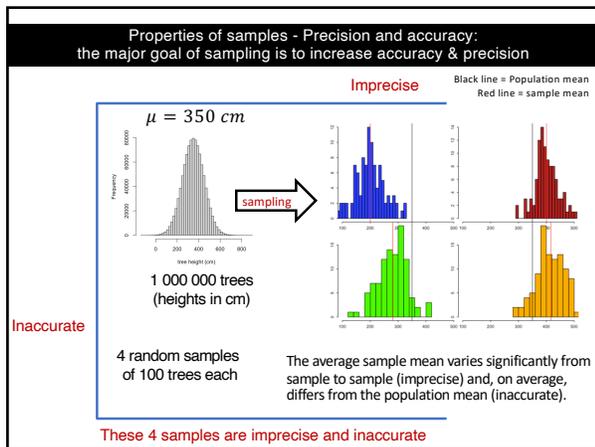
24



25



26



27

Properties of samples - Precision and accuracy:
the major goal of sampling is to increase accuracy & precision

Imprecise

Inaccurate: Sample values (e.g., sample means) tend to be far from the true population value.

Imprecise: Sample values (e.g., sample means) tend to vary widely from each other, regardless of whether they are close to or far from the population value.

Inaccurate

28

Random sampling minimizes bias and makes it possible to measure the amount of sampling error (next lectures)

	Precise	Imprecise
Accurate	<p>Low sampling variation (sampling error) & low bias</p>	<p>High sampling variation (sampling error) & low bias</p>
Inaccurate	<p>Low sampling variation (sampling error) & high bias</p>	<p>High sampling variation (sampling error) & high bias</p>

29

Random sampling minimizes bias and allows for the measurement of sampling error (to be covered in upcoming lectures).

Sample bias: Occurs when certain observational units in the target population have a lower or higher probability of being sampled.

Inferential bias: Arises when the average of all sample values for the statistic of interest (e.g., mean tree height) differs from the true population value. There are many sources of inferential bias, including a lack of random sampling (other sources will be discussed later).

Sampling variation: Refers to the variation that occurs from sample to sample in terms of the statistic of interest. High sampling variation leads to imprecision.

30

Sampling populations - what can go wrong?
 Issues with biased samples based on **sampling of convenience**

Sampling bias occurs when certain members of a population are systematically more likely to be selected in a sample than others

Number of injuries per cat

Number of stories fallen

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

The author proposed that dropping more than 8 floors allows cat to relax and change muscles to cushion their impact

Mehlaff (1987) – Journal of the American Veterinary Medical Association

31

Sampling populations - what can go wrong?
 Issues with biased samples based on **sampling of convenience**

Number of injuries per cat

Number of stories fallen

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

Critics of the study pointed out that instantly fatal falls were not included.

Issue with samples of convenience = existent data not collected for the purposes of the study.

32

Sampling populations - what can go wrong?
 Observational units may vary in other aspects that may lead to sampling biases

Oecologia (2013) 171:339–345
 DOI 10.1007/s00442-012-2426-5

METHODS

Are most samples of animals systematically biased? Consistent individual trait differences bias samples despite random sampling

Peter A. Biro

Sampling bias occurs when certain members of a population are systematically more likely to be selected in a sample than others.

33

Sampling populations - what can go wrong?

Observational units may vary in other aspects that may lead to sampling biases

Volunteer bias

In a large experiment to test the effectiveness of a polio vaccine, schoolchildren were randomly selected to receive either the vaccine or a saline solution as a control.

The vaccine was shown to be effective, but it was later discovered that the rate of polio infection among children in the saline group was higher than that of the general population.

One possible explanation (as suggested by Bland, 2000) is that parents of children who had not been exposed to polio—and therefore had no immunity—were more likely to volunteer their children for the study than parents of children who had already been exposed.

Bland, M. 2000. An Introduction to Medical Statistics, 3rd ed. Oxford, UK: Oxford University Press.

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

34

Sampling populations - what can go wrong?

Observational units may vary in other aspects that may lead to sampling biases

Volunteer bias

In a large experiment to test the effectiveness of a polio vaccine, schoolchildren were randomly selected to receive either the vaccine or a saline solution as a control.

The vaccine was shown to be effective, but it was later discovered that the rate of polio infection among children in the saline group was higher than that of the general population.

One possible explanation (as suggested by Bland, 2000) is that parents of children who had not been exposed to polio—and therefore had no immunity—were more likely to volunteer their children for the study than parents of children who had already been exposed.

Compared to the general population, volunteers may:

- Be more health-conscious and proactive;
- Have lower incomes (especially if volunteers are compensated);
- Be more ill, particularly if the therapy carries risk, as individuals with severe illnesses may be willing to try anything;
- Have more free time (e.g., retirees or unemployed individuals are more likely to participate in telephone surveys);
- Be more upset or angry, as people who are dissatisfied may be more inclined to express their views (e.g., surveys);
- And so on.

Bland, M. 2000. An Introduction to Medical Statistics, 3rd ed. Oxford, UK: Oxford University Press.

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

35

Look into notes & additional material in the WebBook

Survivorship bias: great video explaining sample bias (also covered in Whitlock & Schluter). This is a great video where wrong understanding of sampling can lead to wrong decisions.



Watch on YouTube

36