


A snap demonstration of why numeracy is key to society



In the 1980s, A&W tried to compete with the McDonald's Quarter Pounder by selling a 1/3 pound burger at a lower cost. The product failed, because most customers thought 1/4 pound was bigger.


1

Lecture 8: Estimating with uncertainty, but with a degree of certainty (i.e., with some confidence).

Statistics is the science of aiding decision-making with incomplete information

"While nothing is more uncertain than a single life, nothing is more certain than the average duration of a thousand lives"

Elizur Wright (mathematician & "the father of life insurance")



Statistics is the study of uncertainty

2

Statistics - like life itself - is all about making big conclusions from (small) samples.

One primary goal of statistics is to estimate (infer) an unknown quantity (parameter) of a population based on sample data.

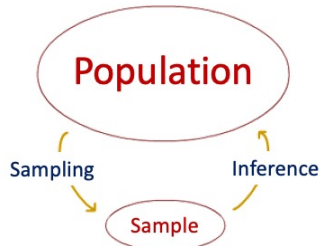
Estimation involves inferring a population parameter (e.g., mean, standard deviation, median) from a sample.

We use estimates to make decisions. Statistics is fundamentally the science of making decisions with incomplete knowledge, often using samples from populations of unknown sizes.

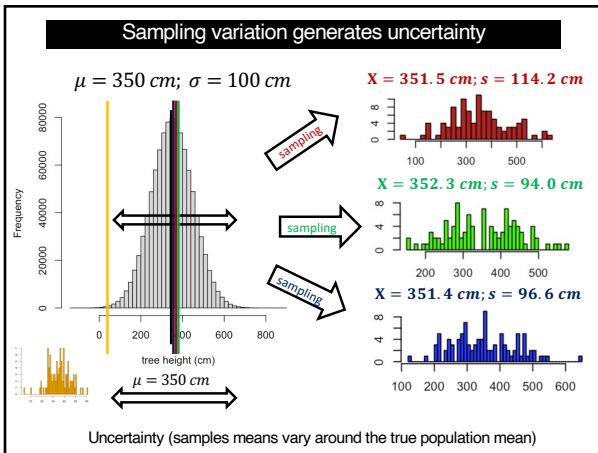
However, sample-based statistics (e.g., mean, median, standard deviation) vary from one sample to another. This variation introduces uncertainty, known as sampling variation.

3

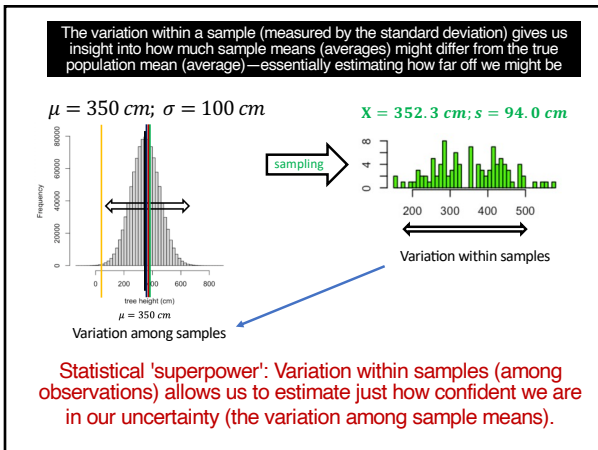
How to estimate with uncertainty, but with some degree of certainty (i.e., with some confidence)?



4



5



6

Population parameters versus sample estimates

A parameter describes a quantity in a statistical population, while an estimate (or statistic) is a similar quantity derived from a sample.

For example, the mean of a population is a parameter, whereas the mean of a sample is an estimate (or statistic) of the population mean.

Similarly, the standard deviation of a population is a parameter, and the standard deviation of a sample is an estimate (or statistic) of the population's standard deviation.

7

Estimating with uncertainty (i.e., error around the true parameter)

An estimate (derived from a sample) is rarely, if ever, exactly the same as the population parameter being estimated—especially in large populations—because sampling is influenced by chance.

For example, two people could sample 100 trees from the same forest and get different mean values. Neither of these sample means will be exactly equal to the population mean.

The critical question in statistics is: **In the face of uncertainty (due to random chance), how much can we trust an estimate and the decisions based on it?** In other words, how accurate is the estimate (i.e., how close is the sample value to the true population value)?

The goal is to deal with uncertainty with a degree of certainty!

8

How to estimate with uncertainty, but with some degree of certainty (i.e., with some confidence)?

We need to understand the properties of estimators (such as the mean, variance, and standard deviation).

These **properties** are examined through the sampling distribution of the statistic or estimate of interest (e.g., sample mean, standard deviation).

A sampling distribution represents the probability distribution of an estimate based on random sampling from the population. It shows what we might observe if we were to repeatedly sample from the population.

While sampling distributions resemble frequency distributions, sampling distributions are made of probabilities instead of frequencies.

9

Statistical symbols

μ = population mean (we say "mu", Greek alphabet).
 σ = population standard deviation (we say "sigma").
 σ^2 = population variance (we say "sigma squared").

10

Important statistical symbols regarding inference

μ = population mean (we say "mu", Greek alphabet).
 σ = population standard deviation (we say "sigma").
 σ^2 = population variance (we say "sigma squared").

\bar{X} = sample mean (we say "X bar", Latin or Roman alphabet).
 s = sample standard deviation.
 s^2 = sample variance.

While μ always represent the mean of the population for any variable you're measuring (e.g., X), the symbol for the sample mean (as discussed before) can vary depending on the variable. For example, it might be written as \bar{X} for the mean of X, or \bar{Y} for the sample mean of Y. However, the key is that it always includes a bar on top of the variable, regardless of which variable you're referring to.

11

**Properties of sampling distributions -
the case of a tiny statistical population of 5 numbers**

1,2,3,4,5; population mean (parameter) = 3.0

All possible 15 samples (with replacement) and their means for $n = 2$:

(1,1) = 1.0	(1,2) = 1.5	(2,3) = 2.5	(3,4) = 3.5	(4,5) = 4.5
(2,2) = 2.0	(1,3) = 2.0	(2,4) = 3.0	(3,5) = 4.0	
(3,3) = 3.0	(1,4) = 2.5	(2,5) = 3.5		
(4,4) = 4.0	(1,5) = 3.0			
(5,5) = 5.0				

Notice that permutations, i.e., (1,2) = (2,1) are not shown but should be considered

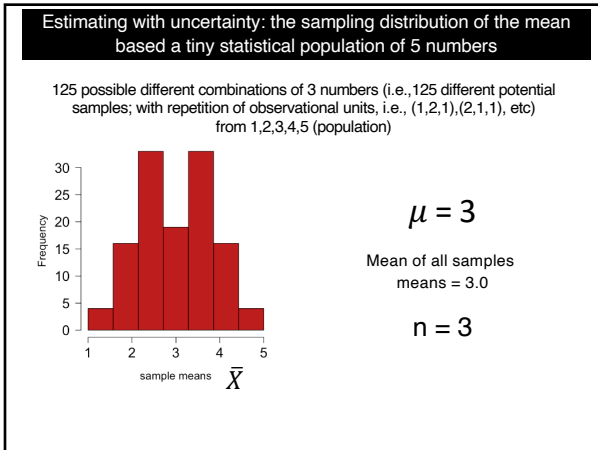
Property 1: The mean of all sample means is always equal to the population mean:

$$(1.0 + 2.0 + 3.0 + 4.0 + 5.0 + 1.5 + 2.0 + 2.5 + 3.0 + 2.5 + 3.0 + 3.5 + 3.5 + 4.0 + 4.5) / 15 = 3.0$$

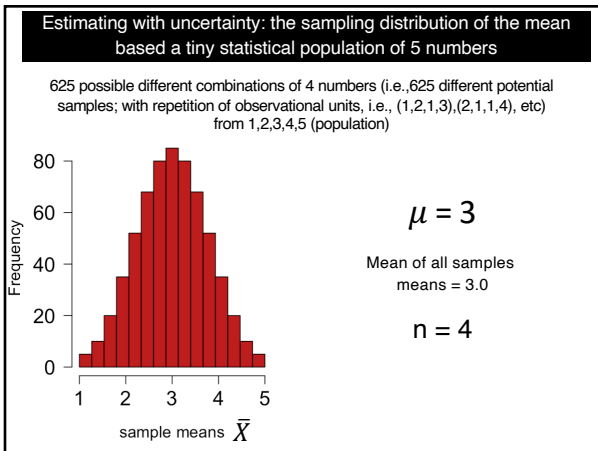
Sample means of the sample population varied from 1.0 to 5.0

sample size (i.e., number of observational units) is represented by the letter "n". Here, $n = 2$ observational units.

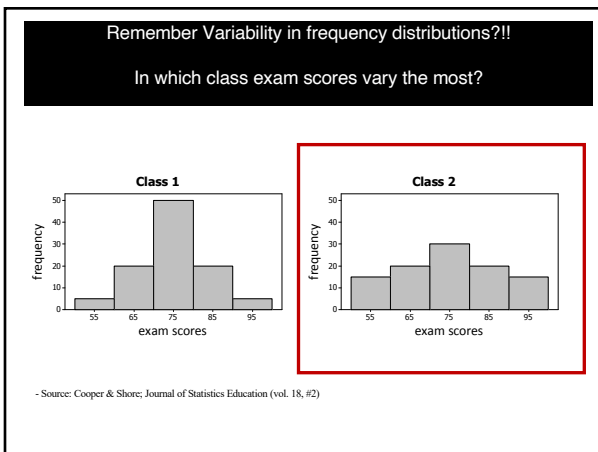
12



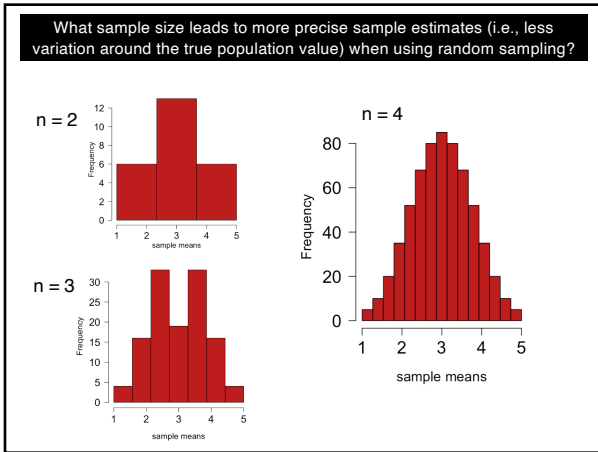
16



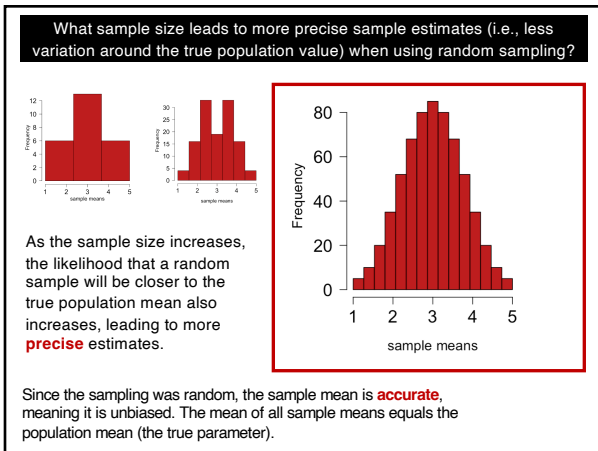
17



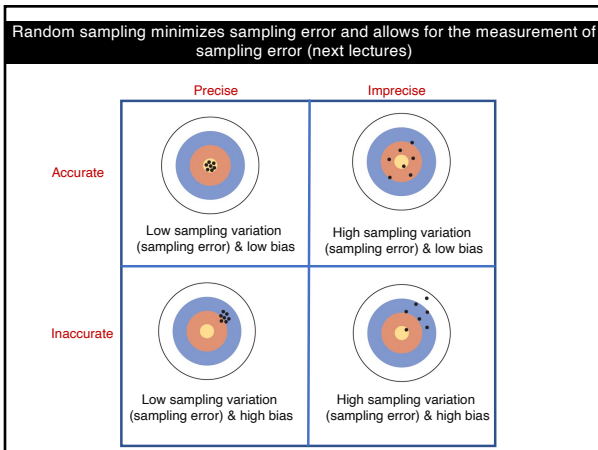
18



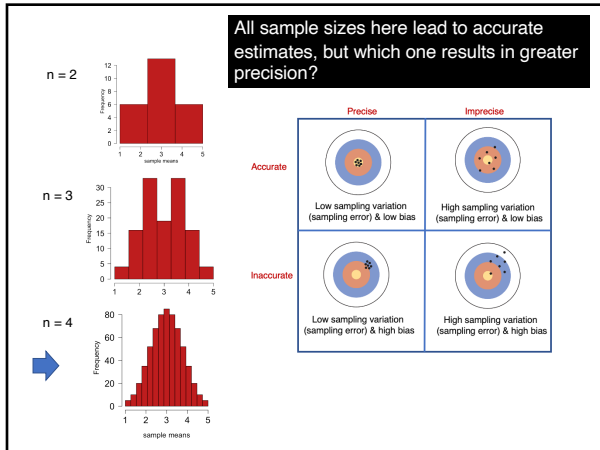
19



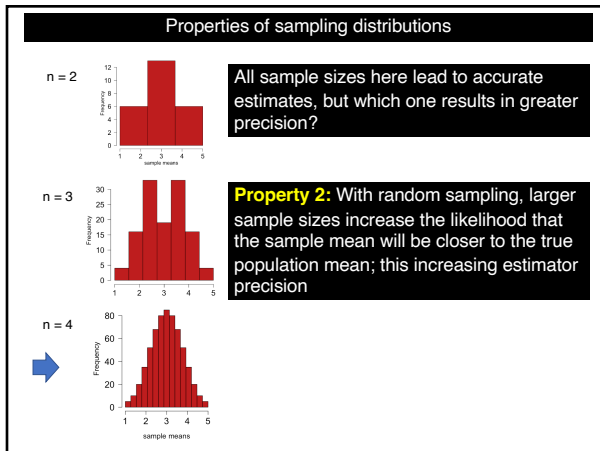
20



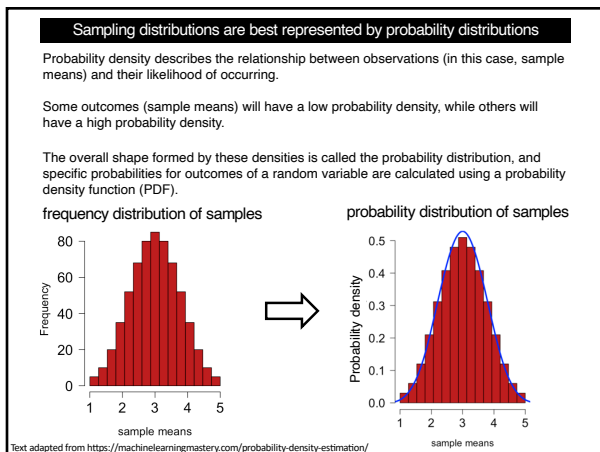
21



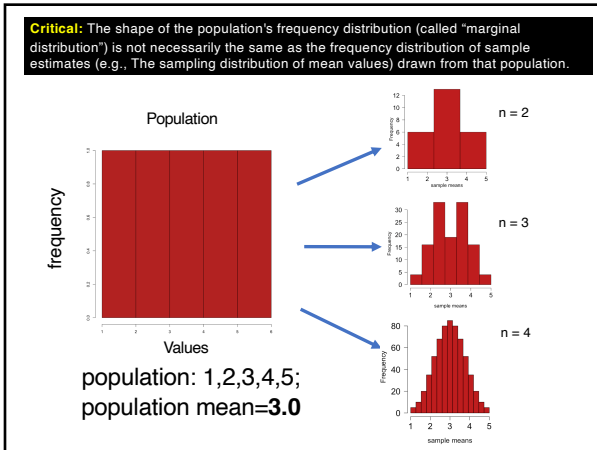
22



23



24



25



26

The length of protein-coding genes in humans is a rare example of an almost complete statistical population in biology

The International Human Genome Project produced the complete DNA sequence for all 23 human chromosomes, each containing millions of nucleotides and more than 23,000 protein-coding genes. The project began in 1990 and was completed in 2006 with the sequencing of the last chromosome. For BIOL 322 tutorials, the available data includes 20,290 genes.

Chromosome
DNA
Genes

27

The length of human genes

It involves the length of almost all human genes, i.e., these is very close to the true population of genes!

Names	Parameter	Value (nucleotides)
Mean (μ)	μ	2622.0
Standard deviation (σ)	σ	2036.9

Frequency distribution of gene lengths in the "known" human genome

28

In real situations, we typically don't know the parameter values of the study population, but in this case, we (almost) do!

So, we'll take advantage of this gene population to illustrate the processes of sampling, uncertainty, accuracy, precision, and how to estimate with uncertainty—yet with some level of confidence!

Names	Parameter	Value (nucleotides)
Mean	μ	2622.0
Standard deviation	σ	2036.9

29

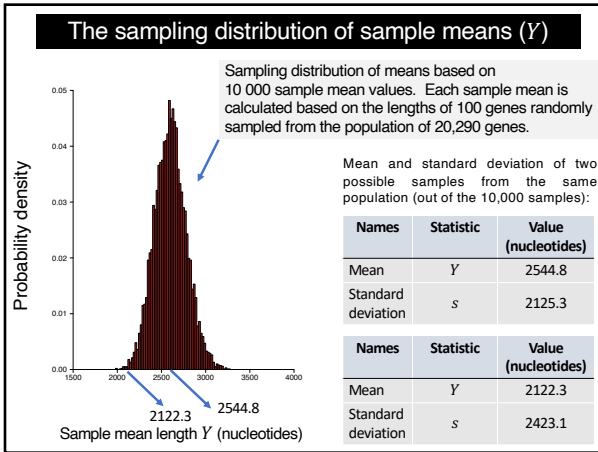
Estimating mean gene length with a random sample of 100 genes (random sampling out of 20,290 genes)

Names	Statistic	Value (nucleotides)
Mean	\bar{Y}	2544.8
Standard deviation	s	2125.3

Frequency distribution of gene lengths in a unique random sample of $n = 100$ genes from the human genome.

Imagine a group in Canada and another in France in 1985 working on the same problem, i.e., estimating the average gene length in the human genome; they would have different sample means

30



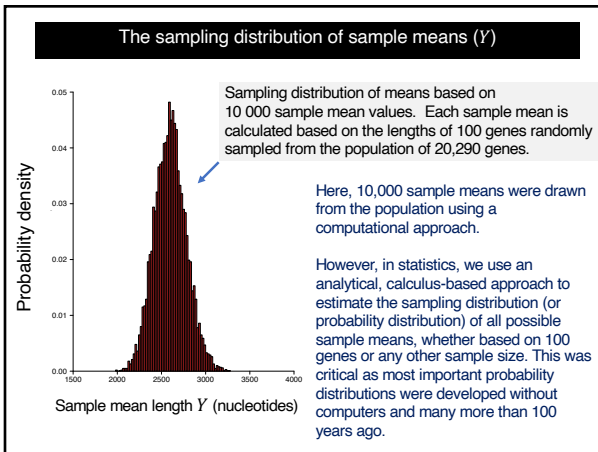
31

Estimating mean gene length with a random sample of 100 genes (random sampling out of 20,290 genes) – variation due to pure chance (i.e., random sampling)

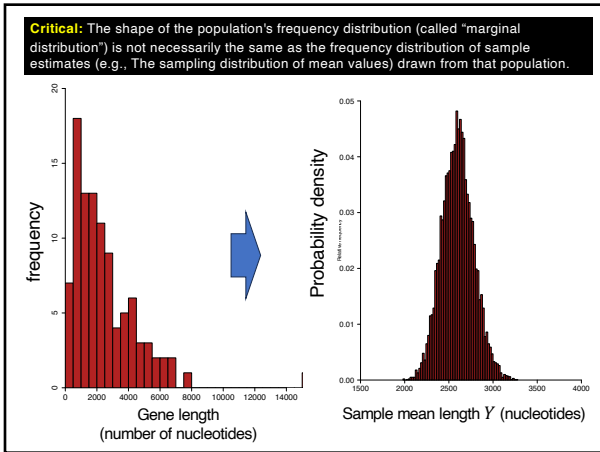
Population			Sample		
Names	Parameter	Value (nucleotides)	Names	Statistic	Value (nucleotides)
Mean	μ	2622.0	Mean	Y	2544.8
Standard deviation	σ	2036.9	Standard deviation	s	2125.3

The sample mean is approximately 77 nucleotides shorter than the true population value. We shouldn't be surprised that the sample estimates differ from the population parameter; such differences are virtually inevitable due to random sampling variation.

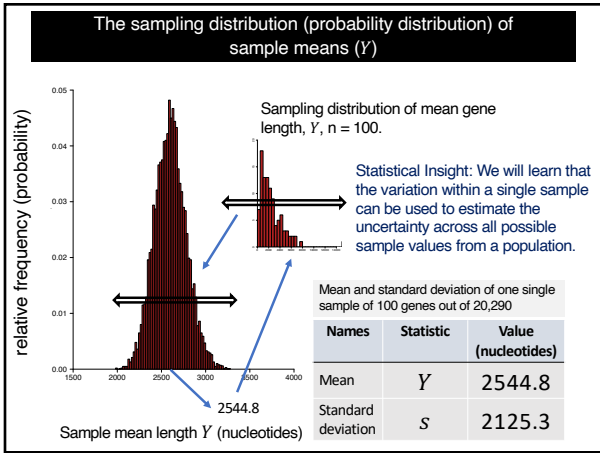
32



33



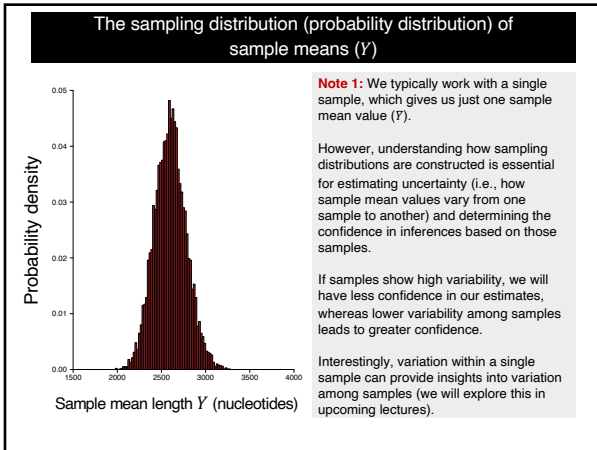
34



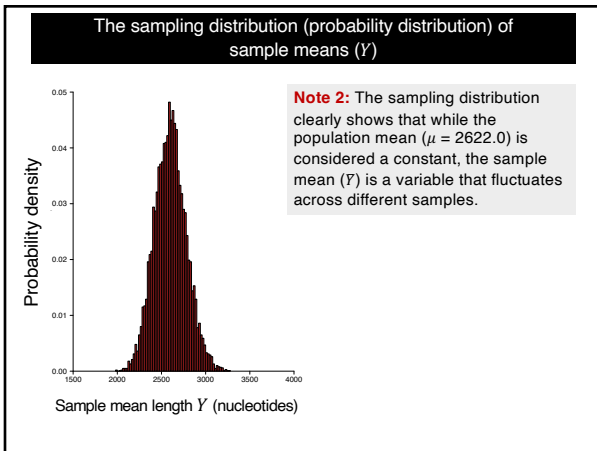
35



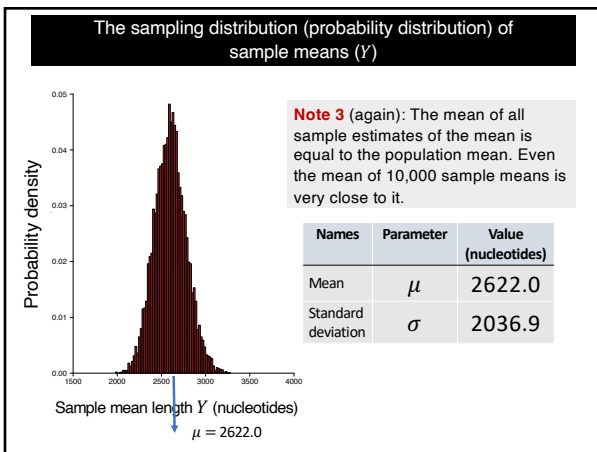
36



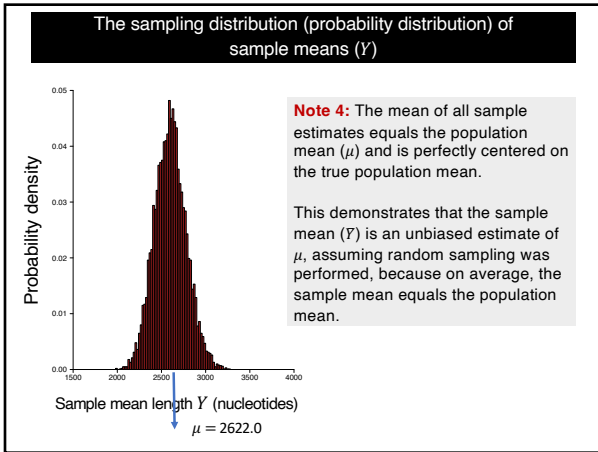
37



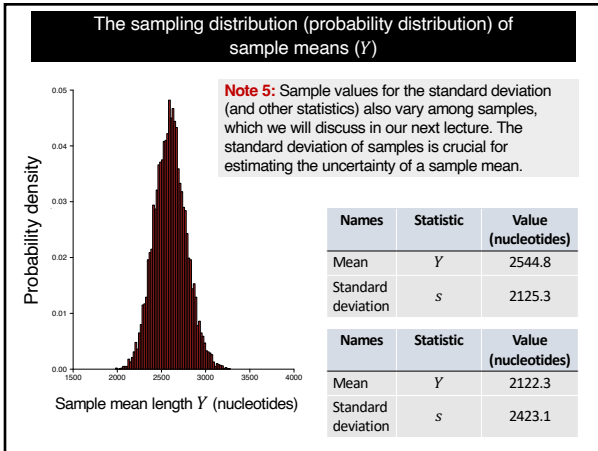
38



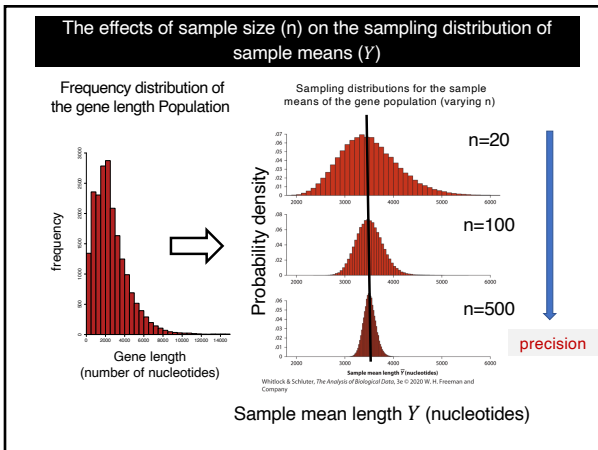
39



40



41



42