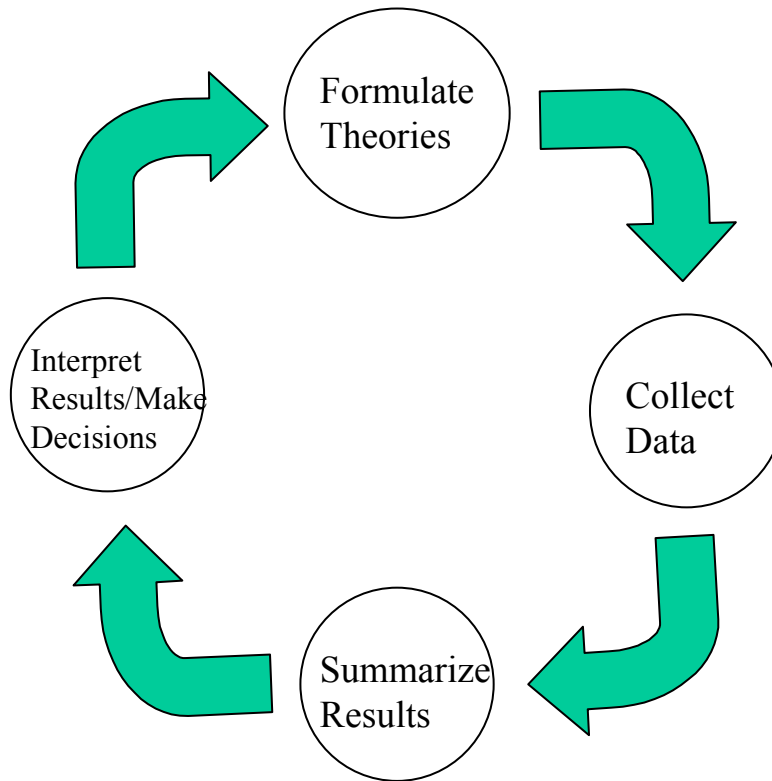


Chapter Goals

To use information from the sample to make inference about the population.



Statistical Sampling

Definition *Parameters are numerical descriptive measures of populations.*

Definition *Statistics are numerical descriptive measures of samples (more generally, a quantity computed from the observations in a sample).*

Definition *Estimators are sample statistics that are used to estimate the unknown population/process parameter. E.g., the sample mean \bar{X} is a point estimator of the population mean μ . Note that \bar{X} is also a random variable.*

Question? How close is our sample statistic to the true, but unknown population parameter?

Notations

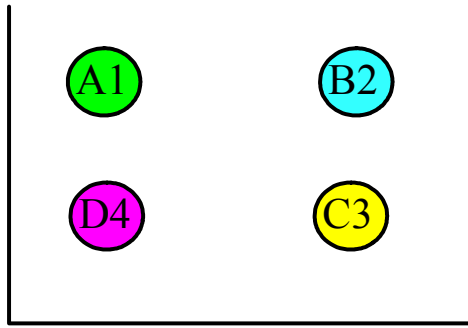
Parameters	Statistics
μ	$\bar{X}, Md, \text{mode}, \hat{\mu}$
σ^2	$s^2, \hat{\sigma}$
p	\hat{p}

Sampling Distributions

Objective: To find out how the sample mean \bar{X} varies from sample to sample. In other words, we want to find out the sampling distribution of the sample mean.

Definition *The sampling distribution of a sample statistic is the distribution of the values of the statistic in all possible samples of the same size n taken from the population.*

Sampling Distribution: An Example



Population Characteristics

Population consists of $N = 4$ items A , B , C , and D , with values 1, 2, 3, and 4 respectively.

Summary measures:

$$\mu = \frac{\sum_{i=1}^4 X_i}{4} = \frac{1+2+3+4}{4} = 2.5$$

$$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N} = \frac{(1-2.5)^2 + (2-2.5)^2 + (3-2.5)^2 + (4-2.5)^2}{4} = 1.25$$

$$\sigma = \sqrt{\sigma^2} = \sqrt{1.25} = 1.118$$

Sampling Distribution of the Sample Mean

Now let us draw all possible samples of size $n = 2$.

16 possible samples				
1 st Observation	2 nd Observation			
	A	B	C	D
A	AA	AB	AC	AD
B	BA	BB	BC	BD
C	CA	CB	CC	CD
D	DA	DB	DC	DD

16 corresponding sample means				
1 st Observation	2 nd Observation			
	A	B	C	D
A	1.0	1.5	2.0	2.5
B	1.5	2.0	2.5	3.0
C	2.0	2.5	3.0	3.5
D	2.5	3.0	3.5	4.0

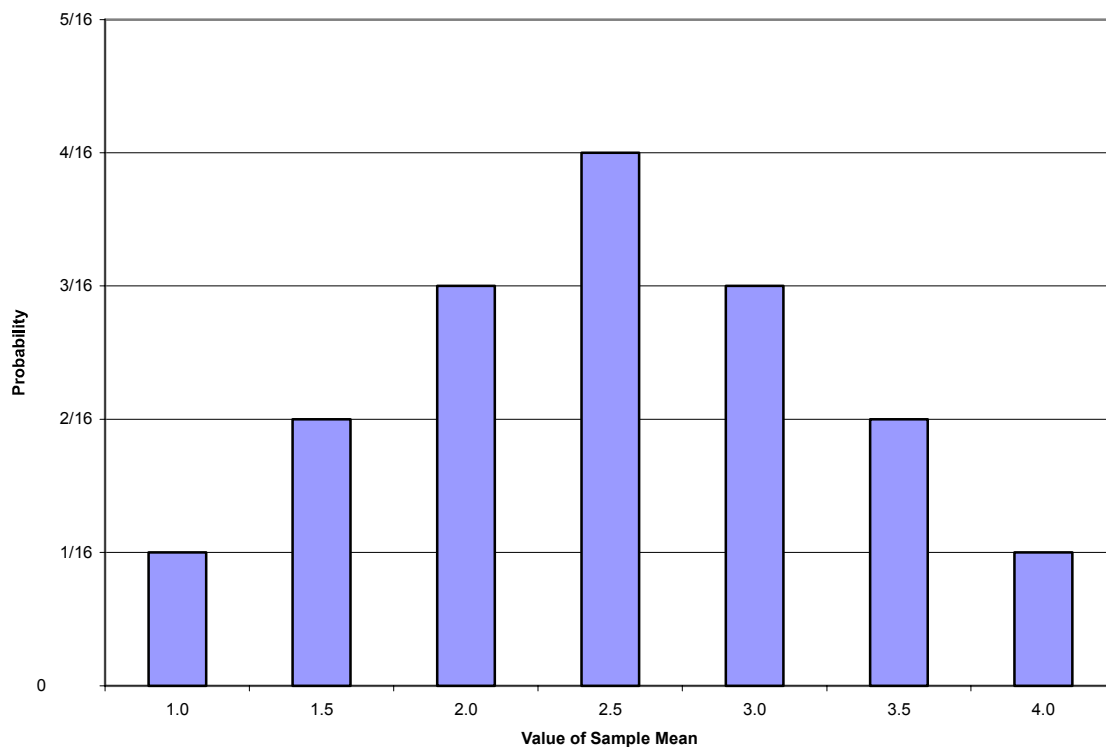
Summary measures:

$$\mu_{\bar{X}} = \frac{\sum_{i=1}^{16} \bar{X}_i}{N} = \frac{1.0+1.5+\dots+4.0}{16} = 2.5$$

$$\sigma_{\bar{X}}^2 = \frac{\sum_{i=1}^{16} (\bar{X}_i - \mu_{\bar{X}})^2}{N} = \frac{(1.0-2.5)^2 + (1.5-2.5)^2 + \dots + (4.0-2.5)^2}{16} = .625$$

$$\sigma_{\bar{X}} = \sqrt{\sigma_{\bar{X}}^2} = \sqrt{.625} = .79$$

Sampling Distribution of \bar{X}



$$E(\bar{X}) =$$

Comparison of Population and Sampling Distribution

	Population	Sampling
μ (Mean)	2.5	2.5
σ^2 (Variance)	1.25	0.625
σ (Standard Deviation)	1.118	0.79

Question: What is the relationship between the variance in the population and sampling distributions?

Empirical Derivation of Sampling Distribution of \bar{X}

- Select a random sample of n observations from a given population or process.
- Compute \bar{X} .
- Repeat steps (1) and (2) a “large” number of times.
- Construct a relative frequency histogram of the resulting set of \bar{X} 's.

Points to Remember about the Sampling Distribution of \bar{X}

1. The mean of the sampling distribution of \bar{X} is the same as the mean of the population/process being sampled from. That is, $E(\bar{X}) = \mu_{\bar{X}} = \mu_X = \mu$.
2. The variance of the sampling distribution of \bar{X} is equal to the variance of the population/process being sampled from divided by the sample size. That is,
$$E\left[(\bar{X} - \mu_{\bar{X}})^2\right] = \sigma_{\bar{X}}^2 = \frac{\sigma_X^2}{n}$$
3. If the original population is normally distributed, then for **any** sample size n the distribution of the sample mean is also normal. That is, if $X \sim N(\mu, \sigma^2)$, then $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$.
4. If the distribution of the original population is not known, but n is sufficiently "large", the distribution of the sample mean is **approximately** normal with mean and variance given as $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$. This result is known as the **central limit theorem (CLT)**.

Example Let X the length of pregnancy be $X \sim N(266, 256)$.

- What is the probability that a randomly selected pregnancy lasts more than 274 days. I.e., what is $P(X > 274)$?

- Suppose we have a random sample of $n = 25$ pregnant women. Is it less likely or more likely (as compared to the above question), that we might observe a sample mean pregnancy length of more than 274 days. I.e., what is $P(\bar{X} > 274)$?

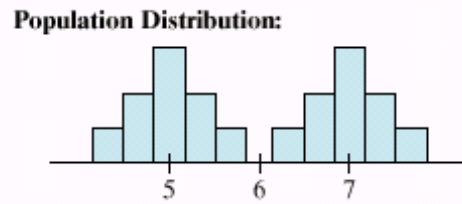
Let's do it! 9.8

Exercise *The model for breaking strength of steel bars is normal with a mean of 260 pounds per square inch and a variance of 400. What is the probability that a randomly selected steel bar will have a breaking strength greater than 250 pounds per square inch?*

Exercise *A shipment of steel bars will be accepted if the mean breaking strength of a random sample of 10 steel bars is greater than 250 pounds per square inch. What is the probability that a shipment will be accepted?*

Let's do it! 9.9

The following histogram shows the population distribution of a variable X . How would the sampling distribution of \bar{X} look, where the mean is calculated from random samples of size 150 from the above population?



Desirable Characteristics of Estimators

Definition An estimator $\hat{\theta}$ is **unbiased** if the mean of its sampling distribution is equal to the population parameter θ to be estimated. That is, $\hat{\theta}$ is an unbiased estimator of θ if $E\{\hat{\theta}\} = \theta$. For example, $E(\bar{X}) = \mu$, therefore \bar{X} is an unbiased estimator of μ .

Definition An estimator is a **consistent** estimator of a population parameter θ if the larger the sample size, the more likely it is that the estimate will be closer to θ . Is \bar{X} a consistent estimator?

Definition The **efficiency** of an unbiased estimator is measured by the variance of its sampling distribution. If two estimators based on the same sample size are both unbiased, the one with the smaller variance is said to have greater relative efficiency than the other.