

Introduction to Probability and Statistics

Three horizontal lines of varying colors (dark blue, green, and blue) are positioned below the title.

Introduction to Probability and Statistics

Three horizontal lines of varying colors (dark blue, green, and blue) are stacked vertically, spanning the width of the slide.

Sampling Distributions

Introduction

- Parameters are numerical descriptive measures for populations.
 - For the normal distribution, the location and shape are described by μ and σ .
 - For a binomial distribution consisting of n trials, the location and shape are determined by p .
- Often the values of parameters that specify the exact form of a distribution are unknown.
- You must rely on the sample to learn about these parameters.

Sampling

Examples:

- A pollster is sure that the responses to his “agree/disagree” question will follow a binomial distribution, but p , the proportion of those who “agree” in the population, is unknown.
- An agronomist believes that the yield per acre of a variety of wheat is approximately normally distributed, but the mean μ and the standard deviation σ of the yields are unknown.
- ✓ If you want the sample to provide reliable information about the population, you must select your sample in a certain way!

Simple Random Sampling

- **The sampling plan or experimental design** determines the amount of information you can extract, and often allows you to measure the **reliability of your inference**.
- **Simple random sampling** is a method of sampling that allows each possible sample of size n an equal probability of being selected.

Example

• There are 89 students in a statistics class. The instructor wants to choose 5 students to form a project group. How should he proceed?



1. Give each student a number from 01 to 89.
2. Choose 5 pairs of random digits from the random number table.
3. If a number between 90 and 00 is chosen, choose another number.
4. The five students with those numbers form the group.

06907	11008	42751	27756	53498
420	69994	98872	31016	
463	07972	18876	20922	
661	10281	17453	18103	
342	53988	53060	59533	
231	33276	70997	79936	
235	03427	49626	69445	
636	92737	88974	33488	
529	85689	48237	52267	
048	08178	77233	13916	

Types of Samples

- Sampling can occur in two types of practical situations:

1. **Observational studies:** The data existed before you decided to study it. Watch out for
 - ✓ **Nonresponse:** Are the responses biased because only opinionated people responded?
 - ✓ **Undercoverage:** Are certain segments of the population systematically excluded?
 - ✓ **Wording bias:** The question may be too complicated or poorly worded.

Types of Samples

- Sampling can occur in two types of practical situations:

2. **Experimentation:** The data are generated by imposing an experimental condition or treatment on the experimental units.

- ✓ **Hypothetical populations** can make random sampling difficult if not impossible.
- ✓ Samples must sometimes be chosen so that the experimenter believes they are **representative** of the whole population.
- ✓ Samples must **behave like random samples!**

Other Sampling Plans

- There are several other sampling plans that still involve **randomization**:

1. **Stratified random sample:** Divide the population into subpopulations or **strata** and select a simple random sample from each strata.
2. **Cluster sample:** Divide the population into subgroups called **clusters**; select a simple random sample of clusters and take a census of every element in the cluster.
3. **1-in-k systematic sample:** Randomly select one of the first k elements in an ordered population, and then select every k -th element thereafter.

Examples



Stratified

- Divide California into counties and take a simple random sample within each county.
- Divide California into counties and take a simple random sample of 10 counties.
- Divide a city into city blocks, choose a simple random sample of 10 city blocks, and interview all who live there.
- Choose an entry at random from the phone book, and select every 50th number thereafter.

Cluster

Cluster

1-in-50 Systematic

Non-Random Sampling Plans

- There are several other sampling plans that do not involve **randomization**. They should **NOT** be used for statistical inference!

1. **Convenience sample:** A sample that can be taken easily without random selection.
 - People walking by on the street
2. **Judgment sample:** The sampler decides who will and won't be included in the sample.
3. **Quota sample:** The makeup of the sample must reflect the makeup of the population on some selected characteristic.
 - Race, ethnic origin, gender, etc.

Sampling Distributions

- Numerical descriptive measures calculated from the sample are called **statistics**.
- Statistics vary from sample to sample and hence are random variables.
- The probability distributions for statistics are called **sampling distributions**.
- In repeated sampling, they tell us what values of the statistics can occur and how often each value occurs.

Sampling Distributions

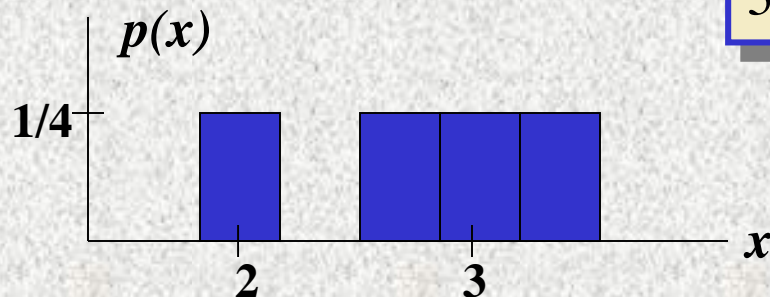
Definition: The sampling distribution of a statistic is the probability distribution for the possible values of the statistic that results when random samples of size n are repeatedly drawn from the population.

Population: 3, 5, 2, 1

Draw samples of size $n = 3$
without replacement

Possible samples	\bar{x}
3, 5, 2	$10/3 = 3.33$
3, 5, 1	$9/3 = 3$
3, 2, 1	$6/3 = 2$
5, 2, 1	$8/3 = 2.67$

Each value of \bar{x} is equally likely, with probability $1/4$



Sampling Distributions

Sampling distributions for statistics can be

- ✓ Approximated with simulation techniques
- ✓ Derived using mathematical theorems
- ✓ The Central Limit Theorem is one such theorem.

Central Limit Theorem: If random samples of n observations are drawn from a nonnormal population with finite μ and standard deviation σ , then, when n is large, the sampling distribution of the sample mean \bar{x} is approximately normally distributed, with mean μ and standard deviation σ / \sqrt{n} . The approximation becomes more accurate as n becomes large.

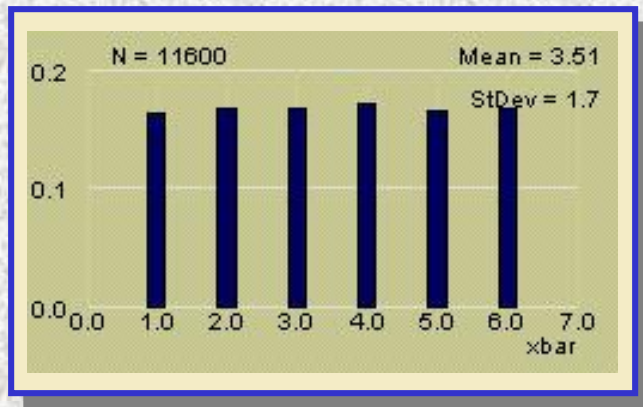


Example

MY

APPLET

Toss a fair coin $n = 1$ time. The distribution of x the number on the upper face is flat or **uniform**.



$$\mu = \sum xp(x)$$

$$= 1\left(\frac{1}{6}\right) + 2\left(\frac{1}{6}\right) + \dots + 6\left(\frac{1}{6}\right) = 3.5$$

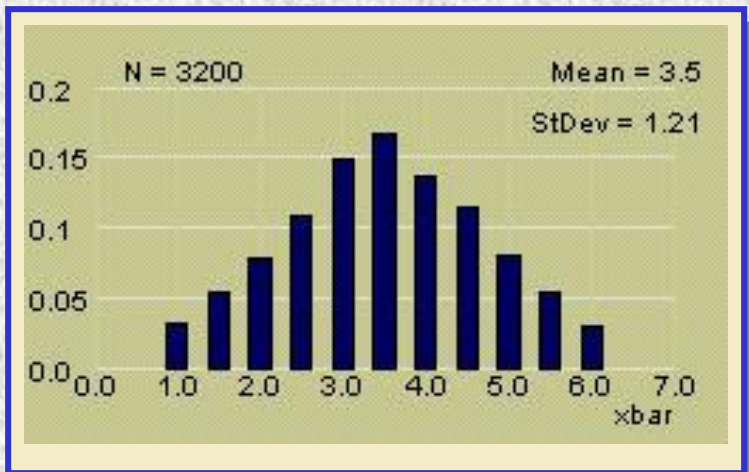
$$\sigma = \sqrt{\sum (x - \mu)^2 p(x)} = 1.71$$



Example



Toss a fair coin $n = 2$ times. The distribution of x the average number on the two upper faces is **mound-shaped**.



$$\text{Mean : } \mu = 3.5$$

Std Dev :

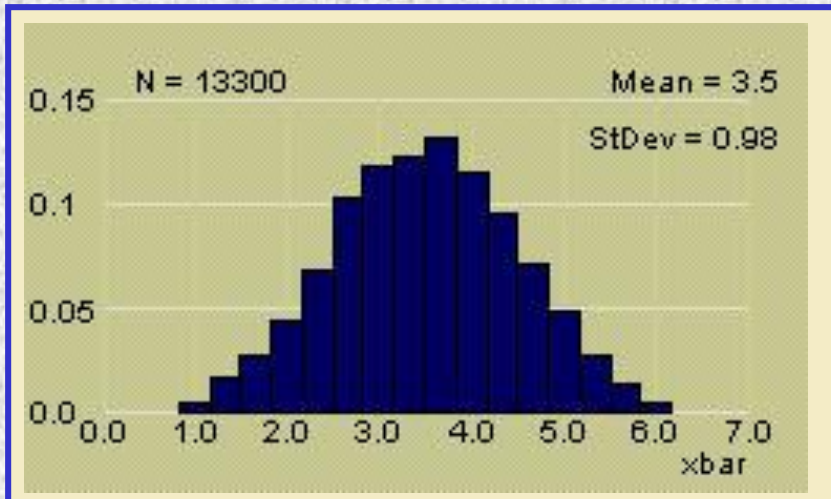
$$\sigma / \sqrt{2} = 1.71 / \sqrt{2} = 1.21$$



Example



Toss a fair coin $n = 3$ times. The distribution of x the average number on the two upper faces is **approximately normal**.

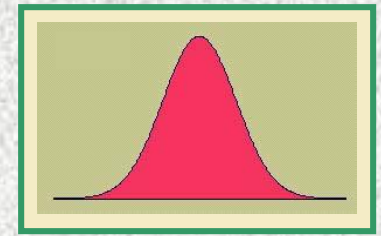


$$\text{Mean : } \mu = 3.5$$

Std Dev :

$$\sigma / \sqrt{3} = 1.71 / \sqrt{3} = .987$$

Why is this Important?



- ✓ The **Central Limit Theorem** also implies that the sum of n measurements is approximately normal with mean $n\mu$ and standard deviation $\sigma\sqrt{n}$.
- ✓ Many statistics that are used for statistical inference are **sums** or **averages** of sample measurements.
- ✓ When n is large, these statistics will have approximately **normal** distributions.
- ✓ This will allow us to describe their behavior and evaluate the **reliability** of our inferences.

How Large is Large?

If the sample is **normal**, then the sampling distribution of \bar{x} will also be normal, no matter what the sample size.

When the sample population is approximately **symmetric**, the distribution becomes approximately normal for relatively small values of n .

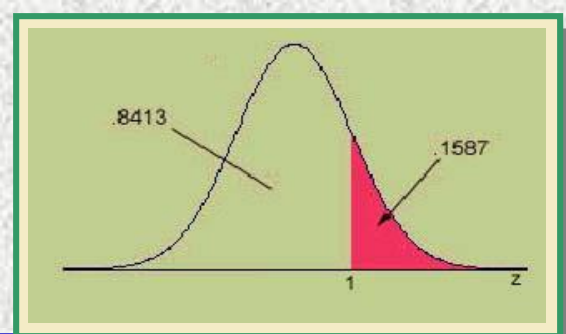
When the sample population is **skewed**, the sample size must be **at least 30** before the sampling distribution of \bar{x} becomes approximately normal.

The Sampling Distribution of the Sample Mean

- ✓ A random sample of size n is selected from a population with mean μ and standard deviation σ .
- ✓ The sampling distribution of the sample mean \bar{x} will have mean μ and standard deviation σ / \sqrt{n} .
- ✓ If the original population is **normal**, the sampling distribution will be normal for any sample size.
- ✓ If the original population is **nonnormal**, the sampling distribution will be normal when n is large.

The standard deviation of \bar{x} is sometimes called the **STANDARD ERROR (SE)**.

Finding Probabilities for the Sample Mean



✓ If the sampling distribution of \bar{x} is normal or approximately normal, *standardize or rescale* the interval of interest in terms of

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

✓ Find the appropriate area using Table 3.

Example: A random sample of size $n = 16$ from a normal distribution with $\mu = 10$ and $\sigma = 8$.

$$\begin{aligned} P(\bar{x} > 12) &= P\left(z > \frac{12 - 10}{8 / \sqrt{16}}\right) \\ &= P(z > 1) = 1 - .8413 = .1587 \end{aligned}$$



Example

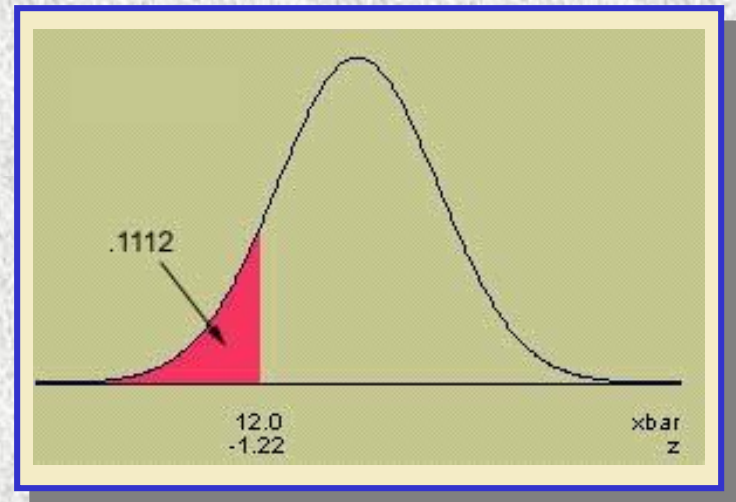


A soda filling machine is supposed to fill cans of soda with 12 fluid ounces. Suppose that the fills are actually normally distributed with a mean of 12.1 oz and a standard deviation of .2 oz. What is the probability that the average fill for a 6-pack of soda is less than 12 oz?

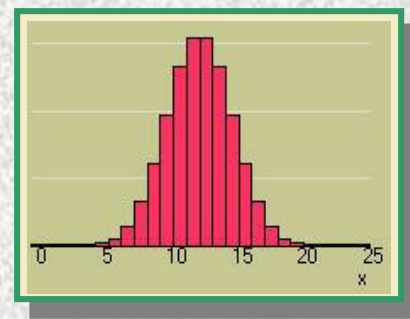
$$P(\bar{x} < 12) =$$

$$P\left(\frac{\bar{x} - \mu}{\sigma / \sqrt{n}} < \frac{12 - 12.1}{.2 / \sqrt{6}}\right) =$$

$$P(z < -1.22) = .1112$$



The Sampling Distribution of the Sample Proportion

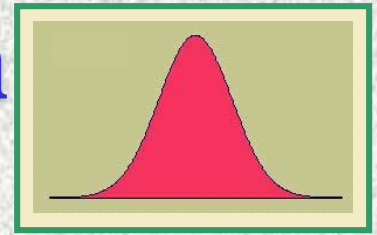


✓ The **Central Limit Theorem** can be used to conclude that the binomial random variable x is approximately normal when n is large, with mean np and standard deviation .

✓ The sample proportion, $\hat{p} = \frac{x}{n}$ is simply a *rescaling* of the binomial random variable x , dividing it by n .

✓ From the Central Limit Theorem, the sampling distribution of \hat{p} will also be **approximately normal**, with a *rescaled* mean and standard deviation.

The Sampling Distribution of the Sample Proportion



✓ A random sample of size n is selected from a binomial population with parameter p .

✓ The sampling distribution of the sample proportion,

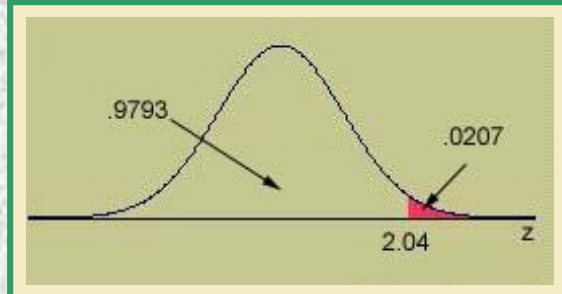
$$\hat{p} = \frac{x}{n}$$

✓ will have mean p and standard deviation $\sqrt{\frac{pq}{n}}$

✓ If n is large, and p is not too close to zero or one, the sampling distribution of \hat{p} will be **approximately normal**.

The standard deviation of p -hat is sometimes called the **STANDARD ERROR (SE)** of p -hat.

Finding Probabilities for the Sample Proportion



✓ If the sampling distribution of \hat{p} is normal or approximately normal, *standardize or rescale* the interval of interest in terms of

$$z = \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}}$$

✓ Find the appropriate area using Table 3.

Example: A random sample of size $n = 100$ from a binomial population with $p = .4$.

$$\begin{aligned} P(\hat{p} > .5) &= P\left(z > \frac{.5 - .4}{\sqrt{\frac{.4(.6)}{100}}}\right) \\ &= P(z > 2.04) = 1 - .9793 = .0207 \end{aligned}$$

Example



The soda bottler in the previous example claims that only 5% of the soda cans are underfilled.

A quality control technician randomly samples 200 cans of soda. What is the probability that more than 10% of the cans are underfilled?

$$n = 200$$

S: underfilled can

$$p = P(S) = .05$$

$$q = .95$$

$$np = 10 \quad nq = 190$$

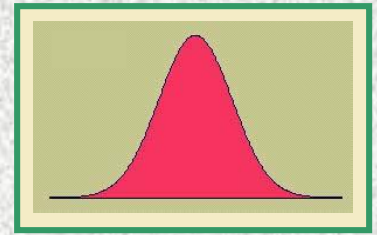
OK to use the normal approximation

$$\begin{aligned} P(\hat{p} > .10) \\ &= P\left(z > \frac{.10 - .05}{\sqrt{\frac{.05(.95)}{200}}}\right) = P(z > 3.24) \\ &= 1 - .9994 = .0006 \end{aligned}$$

This would be very unusual, if indeed $p = .05!$

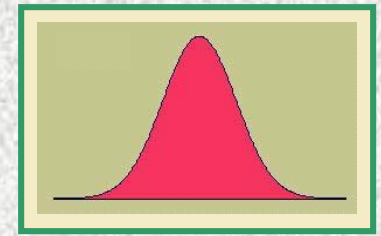
ole

Statistical Process Control



- The cause of a change in the variable is said to be **assignable** if it can be found and corrected.
- Other variation that is not controlled is regarded as **random variation**.
- If the variation in a process variable is solely random, the process is said to be **in control**.
- If out of control, we must reduce the variation and get the measurements of the process variable within specified limits.

The \bar{x} Chart for Process Means

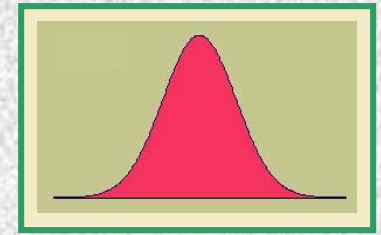


- ✓ At various times during production, we take a sample of size n and calculate the sample mean \bar{x} .
- ✓ According to the CLT, the sampling distribution of \bar{x} should be approximately normal; almost all of the values of \bar{x} should fall into the interval

$$\mu \pm 3 \frac{\sigma}{\sqrt{n}}$$

- ✓ If a value of \bar{x} falls outside of this interval, the process may be out of control.

The \bar{x} Chart

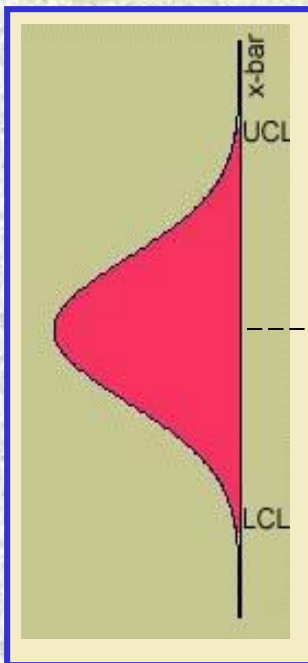


- ✓ To create a control chart, collect data on k samples of size n . Use the sample data to estimate μ and σ .
- ✓ The mean μ is estimated with $\bar{\bar{x}}$, the grand average of all the sample statistics calculated for the nk measurements on the process variable.
- ✓ The standard deviation σ is estimated by s , the standard deviation of the nk measurements.
- ✓ Create the control chart, using a **centerline** and **control limits**.

The \bar{x} Chart

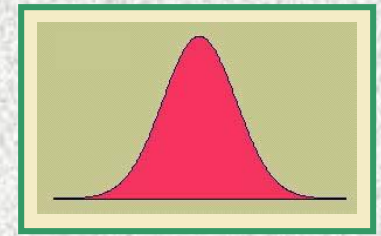
Centerline : $\bar{\bar{x}}$

$$\text{LCL} : \bar{\bar{x}} - 3 \frac{s}{\sqrt{n}} \quad \text{UCL} : \bar{\bar{x}} + 3 \frac{s}{\sqrt{n}}$$



When a sample mean falls outside the control limits, the process may be out of control.

The p Chart for Proportion Defective



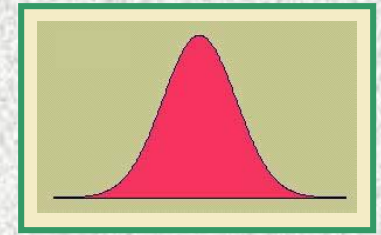
✓ At various times during production, we take a sample of size n and calculate the proportion of defective items, $\hat{p} = x/n$.

✓ According to the CLT, the sampling distribution of \hat{p} should be approximately normal; almost all of the values of \hat{p} should fall into the interval

$$p \pm 3\sqrt{\frac{pq}{n}}$$

✓ If a value of \hat{p} falls outside of this interval, the process may be out of control.

The p Chart



✓To create a control chart, collect data on k samples of size n . Use the sample data to estimate p .

✓The population proportion defective p is estimated with

$$\bar{p} = \frac{\sum \hat{p}_i}{k}$$

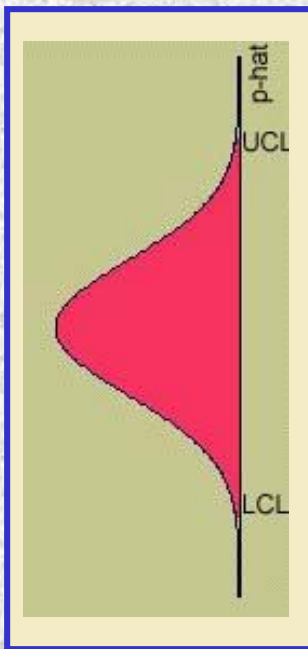
✓the grand average of all the sample proportions calculated for the k samples.

✓Create the control chart, using a **centerline** and **control limits**.

The p Chart

Centerline : \bar{p}

$$\text{LCL} : \bar{p} - 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n}} \quad \text{UCL} : \bar{p} + 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n}}$$



When a sample proportion falls outside the control limits, the process may be out of control.

Key Concepts

I. Sampling Plans and Experimental Designs

1. Simple random sampling

- a. Each possible sample is equally likely to occur.
- b. Use a computer or a table of random numbers.
- c. Problems are nonresponse, undercoverage, and wording bias.

2. Other sampling plans involving randomization

- a. Stratified random sampling
- b. Cluster sampling
- c. Systematic 1-in- k sampling

Key Concepts

3. Nonrandom sampling

- a. Convenience sampling
- b. Judgment sampling
- c. Quota sampling

II. Statistics and Sampling Distributions

1. Sampling distributions describe the possible values of a statistic and how often they occur in repeated sampling.
2. Sampling distributions can be derived mathematically, approximated empirically, or found using statistical theorems.
3. The **Central Limit Theorem** states that sums and averages of measurements from a nonnormal population with finite mean μ and standard deviation σ have approximately normal distributions for large samples of size n .

Key Concepts

III. Sampling Distribution of the Sample Mean

1. When samples of size n are drawn from a normal population with mean μ and variance σ^2 , the sample mean \bar{x} has a normal distribution with mean μ and variance σ^2/n .
2. When samples of size n are drawn from a nonnormal population with mean μ and variance σ^2 , the Central Limit Theorem ensures that the sample mean \bar{x} will have an approximately normal distribution with mean μ and variance σ^2/n when n is large ($n \geq 30$).
3. Probabilities involving the sample mean μ can be calculated by standardizing the value of \bar{x} using

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

Key Concepts

IV. Sampling Distribution of the Sample Proportion

1. When samples of size n are drawn from a binomial population with parameter p , the sample proportion \hat{p} will have an approximately normal distribution with mean p and variance pq/n as long as $np > 5$ and $nq > 5$.
2. Probabilities involving the sample proportion can be calculated by standardizing the value \hat{p} using

$$z = \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}}$$

Key Concepts

V. Statistical Process Control

1. To monitor a quantitative process, use an \bar{x} chart. Select k samples of size n and calculate the overall mean $\bar{\bar{x}}$ and the standard deviation s of all nk measurements. Create upper and lower control limits as
$$\text{LCL: } \bar{\bar{x}} - 3 \frac{s}{\sqrt{n}} \quad \text{UCL: } \bar{\bar{x}} + 3 \frac{s}{\sqrt{n}}$$

If a sample mean exceeds these limits, the process is out of control.

2. To monitor a binomial process, use a p chart. Select k samples of size n and calculate the average of the sample proportions as

$$\bar{p} = \frac{\sum \hat{p}_i}{k}$$

Create upper and lower control limits as

$$\text{LCL: } \bar{p} - 3 \sqrt{\frac{\bar{p}(1-\bar{p})}{n}} \quad \text{UCL: } \bar{p} + 3 \sqrt{\frac{\bar{p}(1-\bar{p})}{n}}$$

If a sample proportion exceeds these limits, the process is out of control.