



Introduction

Probability Distribution
In Chapters 7 and 8, we used knowledge of the population and the parameter(s) of a probability distribution to make probability statements about individual members of the population.

Population & Parameters ⇒ Probability Distribution ⇒Individual

Sampling Distribution
In this chapter, we are going to develop the sampling distribution, wherein knowledge of the parameter(s) and some information about the distribution allow us to make probability statements about a sample statistic.

Population & Parameters >> Sampling Distribution** >> Statistic

Statistical Inference
Starting in Chapter 10, we will assume that most population parameters are unknown. A sample drawn from the population will provide the required statistic, and the sampling distribution of such statistic will enable us to draw inferences about the parameter.

Statistic ⇒ Sampling Distribution ⇒Parameter

Sampling Distribution of the Mean

A sampling distribution is created by sampling a population.

There are two methods to create a sampling distribution:

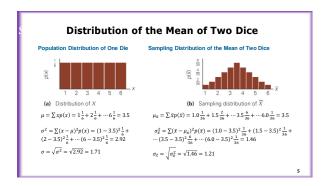
- 1. Theoretical method: apply the rules of probability and the laws of expected value and variance to derive the sampling distribution from sample statistics.
- 2. Empirical method: draw several samples of the same size from a population, calculate the statistic of interest, and then use descriptive techniques to build the sampling distribution.

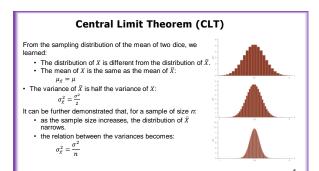
In this chapter we describe the former method and then briefly demonstrate the latter one.

Constructing the sampling distribution of \bar{x}

SOPLE	×	SAMPLE	x	SAMPLE
3.1	18	8.1	10	8.5
1.1	15	3.0	2.5	5.2
3,3	-14	3.1	10	1,1
3,4	13	3.4	13	1.1
1.5	31	1.5	40	3.5
3.6	11	1.6	6.5	3.6
3,1	18	4.1	11	4.1
1.1	20	4.1	3.0	1.1
1.3	13	4.3	13	- 0
2,4	11	4.4	40	6,4
2,1	2.0	4.1	44	- 61
1.4	4.0	4.6	10	5,5







Central Limit Theorem

When the population from which we are selecting a random sample does not have a normal distribution, the central limit theorem is helpful in identifying the shape of the sampling distribution of $\overline{\mathbf{x}}$.

CENTRAL LIMIT THEOREM

In selecting random samples of size n from a population, the sampling distribution of the sample mean \bar{x} can be approximated by a normal distribution as the sample size becomes large.

 ${\bf Sampling\ Distribution\ of\ the\ Sample\ Mean}$

The sampling distribution of the sample mean for large populations can be summarized as:

1.
$$\mu_{\tilde{x}} = \mu$$

2. If the population is infinite:

$$\sigma_{\bar{\chi}} = \frac{\sigma}{\sqrt{n}}$$

3. If the population is finite, then the **finite population correction factor** is applied:

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}$$
 when $N/n < 20$

 If X is normal, then X̄ is normal. If X is nonnormal, then X̄ is approximately normal for sufficiently large sample sizes. The definition of "sufficiently large" depends on the extent of nonnormality of X.

 \rightarrow In many practical situations, a sample size of 30 may be sufficiently large to allow us to use the normal distribution as an approximation for the sampling distribution of \vec{X} .

7

Creating the Sampling Distribution Empirically

We can create the sampling distribution of the mean of two dice empirically by:

- 1. Tossing two fair dice repeatedly.
- 2. Calculating the sample mean for each sample.
- 3. Counting the number of times each value of X occurs.
- 4. Computing the relative frequencies to estimate theoretical probabilities.

Obviously, such empirical approach is not practical because of the excessive amount of time required to toss the dice enough times to make the relative frequencies good approximations for the theoretical probabilities.

Example - Content of a 32-Ounce Bottle

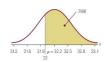
The amount of soda in each 32-ounce bottle is a normally distributed random variable, with mean of 32.2 oz and standard deviation of 0.3 oz.

What is the probability that:

- 1. A single bottle contains more than 32 oz?
- The mean amount of four bottles is greater than 32 oz?

We define:

X = amount of soda per bottle (normal) μ = 32.2 oz, σ = 0.3 oz





. .

Example - Solution

Solution

1. We need to find P(X > 32):

$$\begin{split} P(X > 32) &= P\left(\frac{X - \mu}{\sigma} > \frac{32 - 32.2}{0.3}\right) \\ &= P(Z > -0.67) = 1 - P(-0.67) \\ &= 1 - .2514 = .7486 \end{split}$$

2. We need to find $P(\bar{X} > 32)$ with n = 4:

$$\mu_{\bar{x}} = \mu = 32.2 \text{ and } \sigma_{\bar{x}} = \frac{0.3}{\sqrt{\pi}} = 0.15$$

$$P(\bar{X} > 32) = P\left(\frac{\bar{X} - \mu_{\bar{x}}}{\sigma_{\bar{x}}} > \frac{32 - 32.2}{0.15}\right)$$

$$= P(Z > -1.33) = 1 - P(-1.33)$$

$$= 1 - .0918 = .9082$$

Example: Salaries of Business School Graduates

In the advertisements for a large university, the dean of the School of Business claims that the average salary of the school's graduates one year after graduation is \$800 per week with a standard deviation of \$100.

A student surveys 25 graduates and calculates a mean salary of \$750 per week.

Solution: To verify the dean's claim, we need to calculate $P(\vec{X} < 750)$. The distribution of \vec{X} , the weekly income, is likely to be positively skewed, but not sufficiently so to make the distribution of \vec{X} nonnormal. Thus, we may assume that \vec{X} is no complete the formula of \vec{X} is no complete the \vec{X} in \vec{X} is no complete the \vec{X} in \vec{X} in \vec{X} is no complete the \vec{X} in \vec{X} in \vec{X} in \vec{X} in \vec{X} is no complete the \vec{X} in \vec{X} in \vec{X} in \vec{X} in \vec{X} is no complete the \vec{X} in \vec{X} in \vec{X} in \vec{X} in \vec{X} is no complete the \vec{X} in \vec{X}

$$\begin{split} &\mu_{\widetilde{x}} = \mu = 800, \text{ and } \sigma_{\widetilde{x}} = \sigma/\sqrt{n} = 100/\sqrt{25} = 20 \\ &P(\widetilde{X} < 750) = P\left(\frac{\widetilde{X} - \mu_{\widetilde{X}}}{\sigma_{\widetilde{X}}} < \frac{750 - 800}{20}\right) = P(Z < -2.5) \end{split}$$

The probability is very small. The dean's claim is not justified.



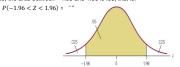
Using the Sampling Distribution for Inference

In preparation for the next chapter on interval estimation, we present **another way** of expressing the probability associated with the sampling distribution.

We define $z_{.025}$ to be the value of Z such that the area to the right of $z_{.025}$ under the standard normal curve is equal to .025. We calculated $z_{.025}$ = 1.96.

Because the standard normal distribution is symmetric about 0, the area to the left of -1.96 is

Thus, the area between -1.96 and 1.96 is .95, that is:



Using the Sampling Distribution for Inference

Because the sampling distribution is normally distributed, we can substitute the equation for ${\it Z}$ into the previous probability statement:

$$P\left(-1.96 < \frac{\bar{X} - \mu_{\bar{X}}}{\sigma / \sqrt{n}} < 1.96\right) = .95$$

Using algebra, we produce:

$$P(\mu_{\bar{X}} - 1.96 \, \sigma / \sqrt{n} < \bar{X} < \mu_{\bar{X}} + 1.96 \, \sigma / \sqrt{n}) = .95$$

If we now plug in mean and standard deviation:

$$P\left(800 - 1.96 \frac{100}{\sqrt{25}} < \bar{X} < 800 + 1.96 \frac{100}{\sqrt{25}}\right) = .95$$

 $P(760.8 < \bar{X} < 839.2) = .95$

We can conclude that there is a 95% probability that a sample mean of size 25 will fall between 5760.8 and \$839.2, which does not include the sample mean of \$750. Thus, the dean's claim is not supported.

Sampling Distribution of a Proportion

Before introduced the binomial distribution, described by the parameter p, the probability of success in any trial.

When p is unknown, we estimate it with the sample proportion (read as p hat): $p = \frac{X}{-}$

Where X is the count of successes and n the sample size. The random variable X is binomially distributed.

However, a binomial random variable is discrete, and it does not lend well to the calculations needed for inference. Next, we show how to approximate a discrete binomial random variable with a continuous normal random variable.

15

Binomial Distribution Example

Consider a binomial distribution with n = 20, and p = .5.

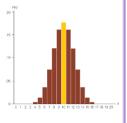
The probability that X = 10, can be written as:

$$P(X = 10) = \frac{n!}{n! (n - X)!} p^{X} (1 - p)^{(n - X)}$$
$$= \frac{20!}{10! (20 - 10)!} p^{10} (1 - p)^{(20 - 10)} = .1762$$

From Section 7-4, we also learned that the binomial distribution parameters are:

$$\mu = np = 20(.5) = 10$$

$$\sigma = \sqrt{np(1 - p)} = \sqrt{20(.5)(.5)} = \sqrt{5}$$



16

Normal Approximation to the Binomial Distribution

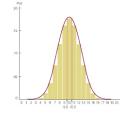
 To approximate the probability that X = 10 using the normal distribution requires that we find the area under the normal curve between 9.5 and 10.5 (0.5 called correction factor).

 $P(X=10) \approx P(9.5 < Y < 10.5)$ where Y is a normal random variable that approximates the binomial random variable X.

We standardize Y to find

 $\begin{array}{l} P(9.5 < Y < 10.5) \\ = P\left(\frac{9.5 - 10}{2.236} < \frac{Y - \mu}{\sigma} < \frac{10.5 - 10}{2.236}\right) \\ = P(-0.2236 < Z < 0.2236) \\ = P(Z < 0.2236) - P(Z < -0.2236) \\ = .5885 - .4115 = .1770 \end{array}$

Which is an excellent approximation.



Omitting the Correction Factor for Continuity

When calculating the probability of individual values of X as we did when we computed the probability that X equals 10 earlier, the **continuity correction factor** must be used.

When computing the probability of a $\ensuremath{\textit{range}}$ of values of X, we can omit the correction factor.

However, the omission of the correction factor will considerably decrease the accuracy of the approximation when the value of X is near the center of the distribution, but it can be tolerated for large sample sizes or when the calculations involve the tails of the

Consider finding $P(X \le 8)$ from the previous example.

Binomial distribution: $P(X \le 8) = .2517$

Normal approximation

with the correction factor: P(Y<8.5) = P(Z<-0.894) = .1857 without the correction factor: P(Y<8) = P(Z<-0.671) = .2511

Approximate Sampling Distribution of a Sample Proportion

Using the laws of expected value and variance, we can determine the mean, variance, and standard deviation of \hat{P} :

Sampling Distribution of a Sample Proportion

 \hat{P} is approximately normally distributed provided that $np \geq 5$ and $n(1-p) \geq 5$

The expected value: $\mathrm{E}(\hat{P})=p$

The variance: $V(\hat{P}) = \sigma_{\hat{P}}^2 = \frac{p(1-p)}{n}$

The standard deviation: $\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$

The standard deviation of \hat{P} is called the **standard error of the proportion**.

Example - Political Survey

In the last election, a state representative received 52% of the votes cast. One year after the election, the representative organized a survey that asked a random sample of 300 people how they would vote in the next election. If we assume that the representative's popularity has not changed, what is the probability that more than half of the sample would vote to re-elect?

Solution

The number of respondents who would vote for the representative is a random binomial variable with n = 300, and p = .52.

We want to find: $P(\hat{P} > .50)$.

 \hat{P} is approximately normal with $E(\hat{P}) = p = .52$ and $\sigma_{\hat{P}} = \sqrt{\frac{p(1-p)}{n}} = \sqrt{\frac{.(52)(48)}{300}} = .0288$

 $P(\hat{P} > .50) = P\left(\frac{\hat{P} - p}{\sigma_{\hat{P}}} > \frac{.50 - .52}{.0288}\right) = P(Z > -0.69) = 1 - P(Z < 0.69) = 1 - .2451 = .7549$

If the level of support has not changed, the probability of a favorable sample exceeds 75%.

20

