

# Sampling Distributions

## Lecture 5

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2     3.5.2     v tibble     3.3.0
v lubridate  1.9.4      v tidyr      1.3.1
v purrr       1.1.0
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

## 1 Outline

- Importance of Sampling Distributions
- Distributions of the Sample Mean
- Distributions of the Difference Between Two Sample Means
- Distributions of the Sample Proportion
- Distributions of the Difference Between Two Proportions

## 2 Importance of Sampling Distributions

### 2.1 Activity

Let us use the “Sampling Words” applet through this [link](#). Suppose we are interested in calculating the average word length of different samples of 10 words from Beyonce’s Crazy in Love.

- Did we get the same average length for all samples?

## 2.2 Sampling Distributions

Repeated samples will yield different values for the statistics.

### **i** Note

We have to view statistics as the sample mean  $\bar{x}$  and sample proportion  $\hat{p}$  as random variables. The probability distribution of these statistics are called sampling distributions. These distributions enable us to:

- Answer probability questions about sample statistics
- Provide the necessary theory for some statistical inference tests.

## 2.3 Sampling Distributions

The sampling distribution of a statistic is the distribution of *all possible values* that can be assumed by some statistic computed from samples of the same size from the population.

### **i** Note

We are usually interested in knowing the functional form (refer to distributions discussed in Chapter 4), mean, and variance.

## 3 Distributions of the Sample Mean

### 3.1 Sampling from Gaussian-distributed populations

For a sample with size  $n$  that comes from a Gaussian-distributed population with mean  $\mu$  and variance  $\sigma^2$ , i.e. the samples  $X_1, X_2, \dots, X_n$  are all independent and identically distributed such that  $X_i \sim N(\mu, \sigma^2)$ , the sample mean can be defined as:

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

The probability distribution of the sample mean, also known as the sampling distribution of the sample mean can be expressed as:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

### 3.2 Implications

#### ! Important

The mean of the sampling distribution,  $\mu_{\bar{X}}$  is equal to the population mean, while the variance,  $\sigma_{\bar{X}}^2$  is reduced by a factor of the sample size  $n$ .

The standard deviation of the sampling distribution can be calculated by taking the square root of the variance of the distribution,  $\sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$ . This is referred to as the **standard error** of the sample mean.

### 3.3 Standardization

The sample mean can also be standardized based on the properties of the sampling distribution. The sample mean  $\bar{X}$  can be transformed into a random variable  $Z$  such that  $Z \sim N(0, 1)$ .

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

In addition,

$$P(\bar{X} \leq \bar{x}) = P(Z \leq \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}) = P(Z \leq z)$$

### 3.4 Sampling from Non-Gaussian Populations

The sampling distribution previously derived could also apply to samples from non-Gaussian populations under certain conditions. These conditions are provided by the **central limit theorem**.

#### i Central Limit Theorem (CLT)

Given a population of any non-Gaussian functional form with a mean  $\mu$  and finite variance  $\sigma^2$ , the sampling distribution of  $\bar{x}$  computed from samples of size  $n$  from this population, will have mean  $\mu$  and variance  $\sigma^2/n$  and will be approximately Gaussian distributed **when the sample size is large**.

#### 💡 Tip

The rule of thumb for sample means is that a sample size of 30 is satisfactory for the central limit theorem, but this maybe too small for skewed distributions. The Gaussian

approximation provided by the CLT becomes better as the sample size increases.

### 3.5 Example

Suppose it is known that in a certain large human population cranial length is approximately Gaussian distributed with a mean of 185.6 mm and a standard deviation of 12.7 mm.

#### 3.5.1 Question

What is the probability that a random sample of size 10 from this population will have a mean greater than 190?

#### 3.5.2 Answer

We can use the sampling distribution because we know the sample is from a Gaussian population. The mean of the sampling distribution,  $\mu_{\bar{X}}$ , should be equal to the population mean. Hence,  $\mu_{\bar{X}} = 185.6$ .

The standard error,  $\sigma_{\bar{X}}$ , can be calculated using the following formula:  $\sigma_{\bar{X}} = \sigma/\sqrt{n} = 12.7/\sqrt{10} = 4.0161$

Using the knowledge that  $\bar{X} \sim N(185.6, 12.7/\sqrt{10})$ , we can use `pnorm` to calculate the probability that the sample mean is greater than 190.

```
pnorm(190,mean=185.6,sd=12.7/sqrt(10),lower.tail=F)
```

```
[1] 0.1366286
```

```
z <- (190-185.6)/(12.7/sqrt(10))  
pnorm(z,0,1,lower.tail=F)
```

```
[1] 0.1366286
```

### 3.6 Example 2

Suppose the hourly number of customer arrivals at a hospital ED has a Poisson distribution with rate parameter  $\lambda = 6$ . Every hour for 48 hours, the number of customer arrivals is counted and recorded.

### 3.6.1 Question

Use the CLT to approximate the probability that the average number of arrivals is between 5 and 8.

### 3.6.2 Answer

The data is non-Gaussian, but we will assume that the CLT holds. Recall that for a Poisson distribution, the mean is equal to the variance. Specifically,  $\mu = 6$ ,  $\sigma^2 = 6$ ,  $n = 48$ . Hence, the mean and the standard error can be expressed as:

$$\mu_{\bar{X}} = 6; \sigma_{\bar{X}} = \sqrt{6/48}$$

The probability that the average number of arrivals in the 48-hour period is between 5 and 8 is:

```
pnorm(8,6,sqrt(6/48)) - pnorm(5,6,sqrt(6/48))
```

```
[1] 0.9976611
```

## 3.7 Exercise

The daily average screen time of elementary school students is assumed to follow a non-Gaussian distribution with 2.6 hours with a standard deviation of 5.3.

### 3.7.1 Question

Use the central limit theorem to calculate the probability that a sample of 250 elementary school students will yield an average screen time between 1.5 hours and 2 hours?

### 3.7.2 Answer

The mean of the sampling distribution is 2.6 hours and the standard error of the mean is  $5.3/\sqrt{250} = 0.3352014$ .

```
pnorm(2,2.6,5.3/sqrt(250)) - pnorm(1.5,2.6,5.3/sqrt(250))
```

```
[1] 0.03621341
```

## 4 Distribution of the Difference Between Two Sample Means

### 4.1 Sample Mean vs. Difference Between Two Sample Means

Consider Populations A and B with the following possible measurements:

A: {0,1,2}; B:{1,2,3}.

If we take a sample with size 2, we can have the following scenario:

Sample A: {0,2}; Sample B: {1,2}

#### Warning

The sample mean of Sample A is 1, and the mean of Sample B is 1.5. However, the difference between the sample means is -0.5, which is not a possible value for the sample means of A and B.

### 4.2 Sampling from Gaussian-distributed populations

When sampling from two Gaussian-distributed populations  $N(\mu_1, \sigma_1^2)$  and  $N(\mu_2, \sigma_2^2)$ , the distribution of the difference between sample means  $\bar{X}_1 - \bar{X}_2$  is a Gaussian distribution with mean  $\mu_1 - \mu_2$  and variance  $\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$ . Mathematically,

$$\bar{X}_1 - \bar{X}_2 \sim N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$$

### 4.3 Standardization

Similar to the sample mean, the difference of the sample means can also be standardized.

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

In addition,

$$P((\bar{X}_1 - \bar{X}_2) \leq (\bar{x}_1 - \bar{x}_2)) = P\left(Z \leq \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}\right)$$

## 4.4 Example

Suppose it has been established that for a certain type of client the average length of a home visit by a public health nurse is 45 minutes with a standard deviation of 15 minutes, and that for a second type of client the average home visit is 30 minutes long with a standard deviation of 20 minutes.

### 4.4.1 Question

If a nurse randomly visits 35 clients from the first and 40 from the second population, what is the probability that the average length of home visit will differ between client type one and client type two by 20 or more minutes?

### 4.4.2 Answer

```
# Calculate the mean and variance of the sampling distribution first.
diff_means <- 45-30
diff_var <- 15^2/35 + 20^2/40

pnorm(20,diff_means,sqrt(diff_var),lower.tail=F)

# OR An alternative solution

z <- (20 - diff_means)/sqrt(diff_var)

pnorm(z,0,1,lower.tail=F)
```

- ①  $\mu_1 - \mu_2$
- ②  $\sigma_1^2/n_1 + \sigma_2^2/n_2$
- ③ Standardization method

```
[1] 0.1086783
```

```
[1] 0.1086783
```

## 4.5 Exercise

Suppose the age of two student organizations were normally distributed. Organization A has an average age of 20.5 and a standard deviation of 3.7, while Organization B has an average of 19.7 and a standard deviation of 4.5.

#### 4.5.1 Question

If 50 students were sampled from each organization, what is the probability that the sample from Organization A is older than Organization B by a value between one and two years?

#### 4.5.2 Answer

```
diff_means <- 20.5-19.7
diff_var <- 4.5^2/50 + 3.7^2/50

pnorm(2,diff_means,sqrt(diff_var))-pnorm(1,diff_means,sqrt(diff_var))
```

```
[1] 0.3314723
```

```
# OR An alternative solution

z1 <- (2 - diff_means)/sqrt(diff_var)
z2 <- (1 - diff_means)/sqrt(diff_var)

pnorm(z1,0,1) - pnorm(z2,0,1)
```

```
[1] 0.3314723
```