

# Difference of Two Proportions

Colby Community College

## Example 1

Consider an experiment for patients who underwent cardiopulmonary resuscitation (CPR) for a heart attack and we subsequently admitted to a hospital.

These patients were randomly divided into a treatment group, where they received a blood thinner or the control group where they did not receive a blood thinner.

The variable of interest was whether they survived for at least 24 hours.

We really have two samples, and hence two sample proportions:

- The proportion of the treatment group that survived:  $\hat{p}_t$
- The proportion of the control group that survived:  $\hat{p}_c$

*How would we determine if blood thinners actually make a difference with these patients?*

We can look at  $\hat{p}_t - \hat{p}_c$ .

But, what we really want to know is, if blood thinners have an effect of heart attack survival rates in the general population?

## Note

The best point estimate for  $p_1 - p_2$  is  $\hat{p}_1 - \hat{p}_2$ .

### Conditions For The Sampling Distribution Of $\hat{p}_1 - \hat{p}_2$ To Be Normal

$\hat{p}_1 - \hat{p}_2$  can be modeled using a normal distribution when:

**Independence:** The data are independent within and between the two groups. Generally this is satisfied if the data come from a randomized experiment.

**Success-Failure:** The success-failure condition holds for both groups, where we check successes and failures in each group separately.

When these conditions are satisfied, the standard error of  $\hat{p}_1 - \hat{p}_2$  is

$$SE = \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}}$$

where  $p_1$  and  $p_2$  represent the population proportions, and  $n_1$  and  $n_2$  represent the sample sizes.

## Confidence Intervals for $\hat{p}_1 - \hat{p}_2$

When the independence and success-failure conditions are met, we can build confidence interval in the same general manner and before:

$$\text{point estimate} \pm z^* \cdot SE$$

$$\Downarrow$$

$$\hat{p}_1 - \hat{p}_2 \pm z^* \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

## Example 2

We can summarize the results from the experiment in Example 1:

	Survived	Died	Total
Control	11	29	50
Treatment	14	26	40
Total	25	65	90

*Is independence satisfied?*

This is a randomized experiment, so yes.

*Are the success-failure conditions satisfied?*

The treatment group had 11 survivals and 29 deaths, and the control group had 14 survivals and 26 deaths. All are more than 10, so yes.

## Example 2 (Continued)

Let us create a 90% confidence interval.

We first need to calculate the point estimate:

$$\hat{p}_t - \hat{p}_c = \frac{14}{40} - \frac{11}{50} = 0.35 - 0.22 = 0.13$$

Next, the standard error:

$$SE \approx \sqrt{\frac{0.35(1 - 0.35)}{40} + \frac{0.22(1 - 0.22)}{50}} = 0.095$$

Recall that the critical value for a 90% confidence is 1.65.

The confidence interval is:

$$\hat{p}_t - \hat{p}_c \pm z^* \cdot SE \rightarrow 0.13 \pm 1.65 \cdot 0.095 \rightarrow (-0.027, 0.287)$$

We are 90% confident that blood thinners have a difference of -2.7 to 28.7 percentage point impact on survival rate for patients.

*What can be conclude about whether blood thinners help or harm?*

Since 0% is in the confidence interval, we don't have enough evidence to say if blood thinners had any impact.

### Example 3

A 5-year experiment was conducted to evaluate the effectiveness of fish oils on reducing cardiovascular events, where each subject was randomized into one of two groups.

We'll consider heart attack outcomes in these patients:

	heart attack	no event	Total
fish oil	145	12788	12933
placebo	200	12738	12938

Let's create a 95% confidence interval.

We first need to calculate the point estimate:

$$\hat{p}_{\text{fish oil}} - \hat{p}_{\text{placebo}} = \frac{145}{12933} - \frac{200}{12938} = 0.0112 - 0.0155 = -0.0043$$

Next, the standard error:

$$SE \approx \sqrt{\frac{0.0112(1 - 0.0112)}{12933} + \frac{0.0155(1 - 0.0155)}{12938}} = 0.00145$$

### Example 3 (Continued)

Recall that the critical value for a 95% confidence is 1.96.

The confidence interval is:

$$\begin{aligned}\hat{p}_{\text{fish oil}} - \hat{p}_{\text{placebo}} \pm z^* \cdot SE &\rightarrow -0.0043 \pm 1.96 \cdot 0.00145 \\ &\rightarrow (-0.0071, -0.0015)\end{aligned}$$

We are 95% confident that fish oils decreases heart attacks by 0.15 to 0.71 percentage points.

*What can we conclude about the effect of fish oils and heart attacks?*

Since the entire interval is negative, we have strong evidence that fish oil supplements reduce heart attacks.



## Example 4

A mammogram is an X-ray procedure used to check for breast cancer. Whether mammograms should be used is part of a controversial discussion.

A 30-year study was conducted with nearly 90,000 female participants. During a 5-year screening period, each woman was randomized to one of two groups:

- The first received regular mammograms screenings.
- The second received regular non-mammogram screenings.

No intervention was made during the following 25 years of the study, and we consider death resulting from breast cancer over the full 30-year period.

	Death from breast cancer?	
	Yes	No
Mammogram	500	44,425
Control	505	44,405

## Example 4 (Continued)

*Is this a experiment or an observational study?*

Because patients were randomized to receive mammograms or a standard breast cancer exam, this is an experiment.

*What are the hypotheses to test whether there was a difference in breast cancer deaths in the mammogram and control groups?*

$H_0$  : The breast cancer death rate for patients screened using mammograms is the same as the breast cancer death rate for patients in the control. ( $p_{\text{mgm}} - p_{\text{ctrl}} = 0$ )

$H_A$  : The breast cancer death rate for patients screened using mammograms is different than the breast cancer death rate for patients in the control ( $p_{\text{mgm}} - p_{\text{ctrl}} \neq 0$ )

## Definition

When the null hypothesis is that  $p_1 - p_2 = 0$ , we use a special proportion called the **pooled proportion** to check the success-failure condition and compute the standard error:

$$\hat{p}_{\text{pooled}} = \frac{\text{number of "successes"}}{\text{number of cases}} = \frac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2}$$

## Example 4 (Continued)

$$\begin{aligned}\hat{p}_{\text{pooled}} &= \frac{\text{\# of patients who died from breast cancer in the entire study}}{\text{\# of patients in the entire study}} \\ &= \frac{500 + 505}{500 + 44,425 + 505 + 44,405} \\ &= 0.0112\end{aligned}$$

## Example 4 (Continued)

Lets check if it's reasonable to model the difference in proportion using a normal distribution in this study?

Because the patients are randomized, they can be treated as independent.

Remember that we check the success-failure condition under the assumption that the null hypothesis is true. We must also check for each group:

$$\begin{aligned}\hat{p}_{\text{pooled}} \cdot n_{\text{mgm}} &= 0.0112 \cdot 44,925 = 503 \geq 10 \\ (1 - \hat{p}_{\text{pooled}}) n_{\text{mgm}} &= 0.9888 \cdot 44,925 = 44,422 \geq 10 \\ \hat{p}_{\text{pooled}} \cdot n_{\text{ctrl}} &= 0.0112 \cdot 44,910 = 503 \geq 10 \\ (1 - \hat{p}_{\text{pooled}}) \cdot n_{\text{ctrl}} &= 0.9888 \cdot 44,910 = 44,407 \geq 10\end{aligned}$$

Since each is at least 10, we can safely model the difference in proportions using a normal distribution.

## Example 4 (Continued)

The point estimate for the difference in breast cancer rates is:

$$\begin{aligned} p_{\text{mgm}} - p_{\text{ctrl}} &= \frac{500}{500 + 44,425} + \frac{505}{505 + 44,405} \\ &= 0.01113 - 0.01125 = -0.00012 \end{aligned}$$

The standard error is:

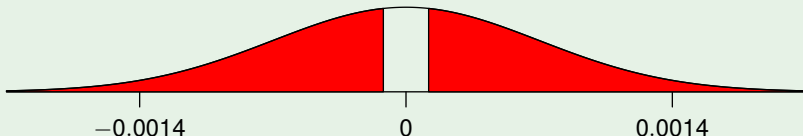
$$\begin{aligned} SE &= \sqrt{\frac{\hat{p}_{\text{pooled}}(1 - \hat{p}_{\text{pooled}})}{n_{\text{mgm}}} + \frac{\hat{p}_{\text{pooled}}(1 - \hat{p}_{\text{pooled}})}{n_{\text{ctrl}}}} \\ &= \sqrt{\frac{0.0112(1 - 0.0112)}{500 + 44,425} + \frac{0.0112(1 - 0.0112)}{505 + 44,405}} = 0.00070 \end{aligned}$$

## Example 4 (Continued)

Next, we calculate the z-score:

$$\begin{aligned} z &= \frac{\text{point estimate} - \text{null value}}{\text{SE}} \\ &= \frac{-0.00012 - 0}{0.00070} = -0.17 \end{aligned}$$

The  $p$ -value is the area in both tails:



The left tail has area 0.4325 and so the total area is 0.8560. We do not reject the null hypothesis.

That is, the difference between breast cancer death rates is reasonably explained by random chance, and we do not observe benefits or harm from mammograms relative to a regular breast exam.

## Example 4 (Continued)

When reviewing this study, or any other study, it's important to keep the following considerations in mind:

- We do not accept null hypothesis, which means we don't have sufficient evidence to conclude that mammograms reduce or increase breast cancer deaths.
- If mammograms are helpful or harmful, the data suggests the effect isn't very large.
- Are mammograms more or less expensive than a non-mammogram breast exam?
  - If one option is much more expensive than the other and doesn't offer clear benefits, then we should lean towards the less expensive option.

## Example 4 (Continued)

- The study's authors also found that mammograms led to overdiagnosis of breast cancer:
  - This means that some breast cancers were found, or thought to be found, but that these cancers would not cause symptoms during the patients' lifetimes.
  - That is, something else would kill the patient before breast cancer symptoms appeared.
  - This means some patients may have been treated for breast cancer unnecessarily, and this treatment is another cost to consider.
  - Overdiagnosis can cause unnecessary physical and emotional harm to patients.