

### 12.3: Inferences for Two Population Proportions

Often, instead of comparing a sample's proportion to a benchmark value, we want to compare the proportions of two different samples, drawn from two different populations.

If the sample sizes are sufficiently large, we can assume that the difference between the two sample proportions,  $\hat{p}_1 - \hat{p}_2$ , is normally distributed.

Suppose we have two samples, Sample A and Sample B, of sizes  $n_1$  and  $n_2$ , respectively.

In Sample A,  $x_1$  is the number of the  $n_1$  observations that possess the characteristic of interest.

In Sample B,  $x_2$  is the number of the  $n_2$  observations that possess the characteristic of interest.

Then the sample proportions are  $\hat{p}_1 = \frac{x_1}{n_1}$  and  $\hat{p}_2 = \frac{x_2}{n_2}$ .

We'll follow this rule of thumb: For sample sizes to be sufficiently large to assume the difference in proportions is normally distributed, all of the quantities  $x_1$ ,  $n_1 - x_1$ ,  $x_2$ ,  $n_2 - x_2$  should be at least 5.

Another way to say this: All four cells (groups) in a contingency table must contain at least 5 data points.

**Example 1:** Suppose we are comparing the number of successful treatments for two groups. Treatment A is given to 20 patients, and is successful in 16 patients. Treatment B is given to 18 patients, and is successful in 12 patients. Can we assume a normal distribution for the difference between proportions?

	Treatment A	Treatment B
Successful	16	12
Unsuccessful	4	6

This cell is less than 5

No, we cannot assume the difference between sample proportions is normally distributed. (So cannot use this procedure)  
(To use this procedure, all cells must be at least 5)

### Sampling distribution for the difference between two sample proportions:

Suppose independent samples of sizes  $n_1$  and  $n_2$  are taken from two populations, in which the proportions of observations possessing a characteristic are  $p_1$  and  $p_2$ .

Then the mean and standard deviation of the difference between sample proportions are:

$$\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$$

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \quad (\text{standard error})$$

Also, the shape of the sampling distribution of  $\hat{p}_1 - \hat{p}_2$  is approximately normal, provided that the sample sizes are sufficiently large.

Rule of thumb: to assume the difference between sample proportion is normally distributed, all of  $x_1$ ,  $n_1 - x_1$ ,  $x_2$ ,  $n_2 - x_2$  should be at least 5.

However, we will almost never know the population proportions  $p_1$  and  $p_2$ , so we cannot calculate the above mean and standard deviation.

Instead, we use statistics from the samples to calculate point estimates for  $\mu_{\hat{p}_1 - \hat{p}_2}$  and  $\sigma_{\hat{p}_1 - \hat{p}_2}$ .

Point estimates for the mean and standard deviation of the difference between sample proportions:

$\hat{p}_1 - \hat{p}_2$  is an unbiased estimate for  $\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2$ .

To calculate a point estimate for  $\sigma_{\hat{p}_1 - \hat{p}_2}$ , we first calculate the pooled sample proportion  $\hat{p}_p$ :

$$\hat{p}_p = \frac{x_1 + x_2}{n_1 + n_2},$$

where  $n_1$  and  $n_2$  are the sample sizes, and  $x_1$  and  $x_2$ , respectively, are the numbers of observations possessing the characteristic of interest. Then the standard error is approximately

$$\sigma_{\hat{p}_1 - \hat{p}_2} \approx \sqrt{\hat{p}_p(1-\hat{p}_p)} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = \sqrt{\hat{p}_p(1-\hat{p}_p) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

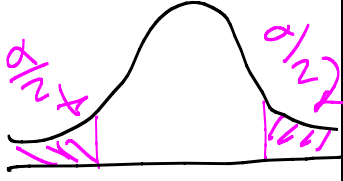

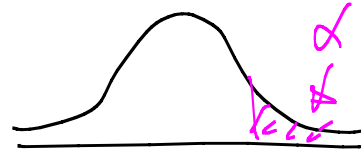
### Hypothesis Testing for Two Population Proportions:

Step 1: Determine the significance level  $\alpha$ .

Step 2: Check that the assumptions are satisfied.

- Simple random sample
- Samples are independent (observations from one sample are not paired with observations from the other sample).
- $x_1$ ,  $n_1 - x_1$ ,  $x_2$ ,  $n_2 - x_2$  are all at least 5.

Step 3: Determine the null and alternative hypotheses.

Two-Tailed Test (most common)	Left-Tailed Test (rare)	Right-Tailed Test (rare)
$H_0 : p_1 = p_2$ $H_a : p_1 \neq p_2$	$H_0 : p_1 = p_2$ $H_a : p_1 < p_2$	$H_0 : p_1 = p_2$ $H_a : p_1 > p_2$
 Rejection Region	 Rejection Region	 Rejection Region

Note: One tailed tests assume that the scenario not listed ( $p_1 < p_2$  for a left-tailed test or  $p_1 > p_2$  for a right-tailed test) is not possible or is of zero interest.

Step 4: Using your  $\alpha$  level and hypotheses, sketch the rejection region.

Step 5: Use a normal curve table to determine the critical value for z associated with your rejection region.

Step 6: Compute the test statistic:

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}_p(1 - \hat{p}_p)} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Step 7: Determine whether the value of  $z$  calculated from your sample (in Step 6) is in the rejection region.

- If  $z$  is in the rejection region, reject the null hypothesis.
- If  $z$  is not in the rejection region, do not reject the null hypothesis.

Step 8: State your conclusion and interpret the results of the hypothesis test.

At least 5 in each cell of table? Yes

**Example 2:** Suppose that a clinical trial for two different cancer drugs is conducted. For drug A, 637 of 2095 patients were cured. For drug B, 702 of 2119 patients were cured. Does this trial provide evidence that Drug B cures a higher percentage of patients than Drug A? Use a 5% level of significance.

$$H_0: p_1 = p_2$$

$$H_a: p_1 < p_2$$

$$\alpha = 0.05$$

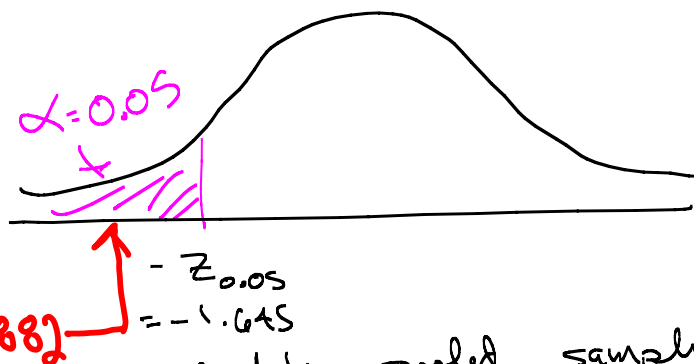
Sample info:

$$\text{Drug A: } \hat{p}_1 = \frac{x_1}{n_1} = \frac{637}{2095} \approx 0.304$$

$$n_1 = 2095$$

$$\text{Drug B: } \hat{p}_2 = \frac{x_2}{n_2} = \frac{702}{2119} \approx 0.331$$

$$n_2 = 2119$$



$$Z = -1.882$$

Calculate pooled sample proportion:

$$\hat{p}_p = \frac{x_1 + x_2}{n_1 + n_2} = \frac{\text{total successes}}{\text{total patients}}$$

$$= \frac{637 + 702}{2095 + 2119} \approx 0.31775$$

$$\hat{q}_p = 1 - \hat{p}_p = 1 - 0.31775 = 0.68225$$

$$\text{standard error: } \sigma_{\hat{p}_1 - \hat{p}_2} \approx \sqrt{\hat{p}_p (1 - \hat{p}_p) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$= \sqrt{0.31775 (0.68225) \left( \frac{1}{2095} + \frac{1}{2119} \right)}$$

$$\approx \sqrt{2.05789 \times 10^{-4}} \approx 0.014345 \text{ (store in calculator)}$$

Calculate test statistic

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{0.014345} \approx -1.882$$

This falls in the rejection region, so we

Reject  $H_0$ .

This sample provides evidence that Drug B cures a higher proportion of patients than Drug A

In the above problem, what would be the meaning of a Type I error?

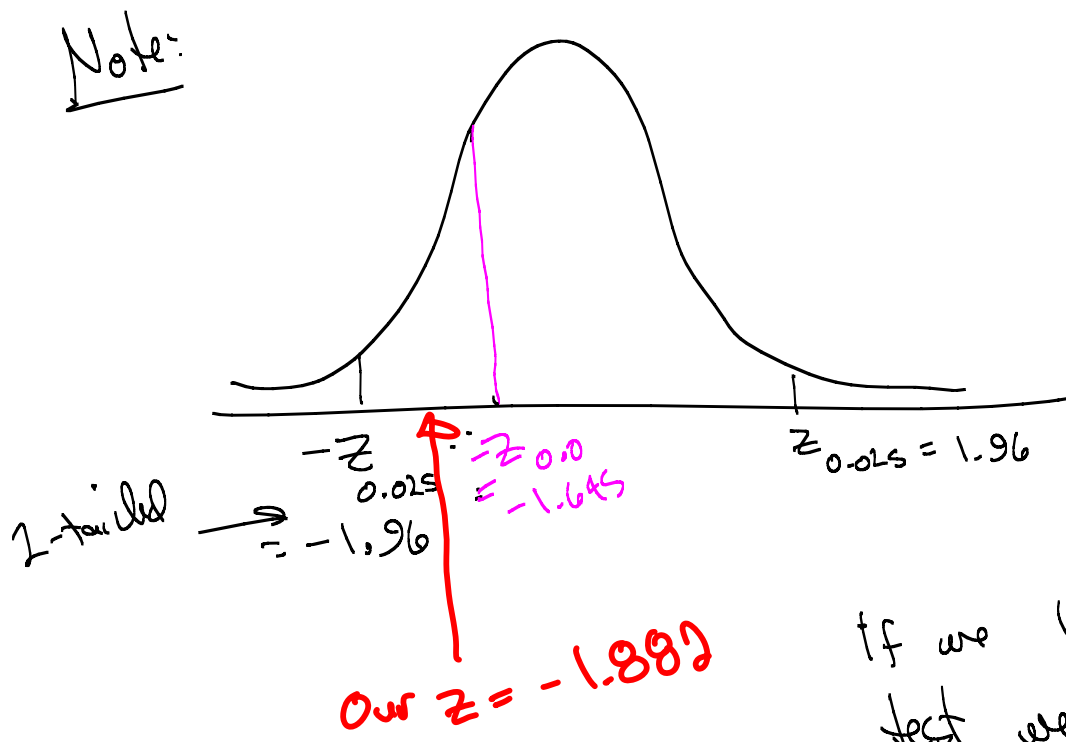
We found a difference in favor of Drug B, but in reality the drugs are equally effective.

(we reject  $H_0$  when  $H_0$  is true.)

A Type II error?

we say the drugs are equally effective when in reality Drug B is better.

(we don't reject  $H_0$  when we should have rejected  $H_0$  ... when  $H_0$  is false)



If we had done a 2-tailed test, we would not have rejected.