

# One-sample test of proportions

## **The Setting:**

Individuals in some population can be classified into one of two categories. You want to make inference about the proportion in each category, so you draw a sample.

## **Examples:**

“Drug X is administered to 100 patients with a particular disease. 50 improve. Test whether this drug is better than drug Y, which is known to produce improvement in 45% of patients.”

“In a poll of 500 voters, 200 say that they will support a particular candidate in the election. Give a 95% confidence interval for the proportion of all voters who will support the candidate.”

# One-sample test of proportions

Inferential Procedures:

- Hypothesis Test
  - critical value method
  - p-value method
  - confidence interval method
- Confidence Interval

# One-sample test of proportions

Hypothesis Test:

Specify the null and alternative

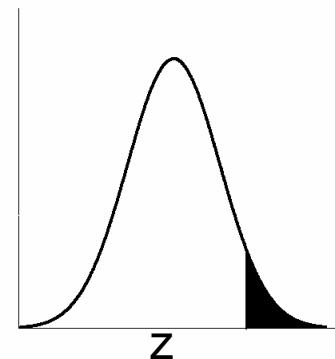
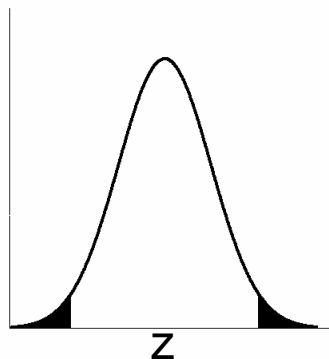
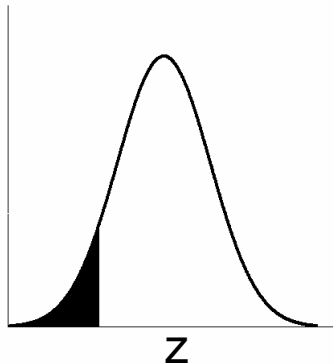
Define rejection region for the test statistic  $z$

$$H_0: p = p_0$$

$$H_A: p < p_0$$

$$H_A: p \neq p_0$$

$$H_A: p > p_0$$



# One-sample test of proportions

Test Statistic: 
$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1 - p_0)}{n}}}$$

If  $H_0$  is true,  $z$  comes from standard normal distribution

## Critical value method:

Choose  $z^*$  from normal table to define rejection region

If  $z$  falls in rejection region, reject  $H_0$

## P-value method:

Use the normal table to determine probability of a value at least as extreme as  $z$  coming from standard normal distribution. Report this p-value and your conclusion.

# One-sample test of proportions

Level-C Confidence Interval:

$$p = \hat{p} \pm z * \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

To perform a two-sided test of  $H_0: p = p_0$  at significance level  $\alpha$ , you can calculate the level- $(1 - \alpha)$  confidence interval and check whether  $p_0$  falls within the interval.

# Two-sample tests: overview

Types of two-sample tests:

- two-sample z-test
- two-sample t-test
- two-sample test of proportions

# Two-sample z-test

## **The Setting:**

You want to determine whether two different populations have the same mean for some variable of interest, so you draw two independent samples. Furthermore, you know the variances in these two populations *without error*.

In real life, this would only tend to happen if your samples were so large (e.g., 5,000) that you could pretend that the sample variances were really the true variances

## **Example:**

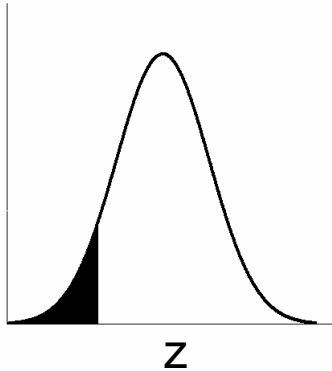
“SAT scores for students who use a particular prep course have an SD of 100, while those for students who do not use the course have an SD of 150. You sample 50 students who take the course and 50 who do not. Test whether taking the course changes the mean score for students”

# Two-sample z-test

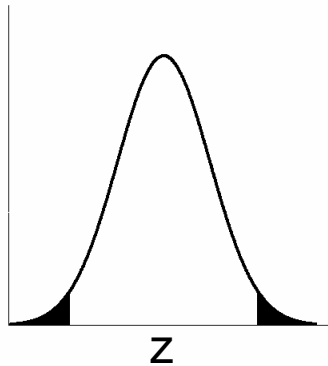
Hypothesis Test

$$H_0: \mu_x = \mu_y$$

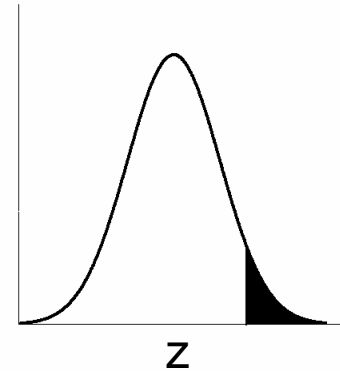
$$H_A: \mu_x < \mu_y$$



$$H_A: \mu_x \neq \mu_y$$



$$H_A: \mu_x > \mu_y$$



Where the test statistic is

$$z = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}}$$



# Two-sample t-test

## **The Setting:**

You want to determine whether two different populations have the same mean value for some variable of interest. You draw two independent samples to test this.

Note that in general we will use this test instead of the z-test, since we will not know the population standard deviations without error.

## **Example:**

“You want to determine whether students who take a given prep course have, on average, the same SAT score as those who do not. You sample 50 students who take the course and 50 who do not and record their SAT scores”

# Two-sample t-test

## Hypothesis Test

- Null and alternative hypotheses same as for two-sample z-test

- Test statistic is now 
$$t = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}}$$

- If  $H_0$  is true,  $t$  follows a t-distribution with  
 $df = \min(n_x - 1, n_y - 1)$

# Two-sample t-test

- **Common Mistake:** Sometimes you will have what appear to be two independent samples, but are really two observations each on a single set of units. Use one-sample methods for inference on such data. This sample will consist of the differences between the two values for each subject.

# Not a Two-sample t-test

## Example:

“You want to determine whether a drug lowers patient cholesterol. You record the cholesterol levels of 100 volunteers, then administer the drug for three months. After this period you again measure the patients’ cholesterol values”

## Strategy:

Define  $D$  = change in a patient’s cholesterol after taking drug.

Now do a one-sample t-test:

$$H_0: \mu_D = 0$$

$$H_A: \mu_D < 0$$

The test statistic will then be  $t = \frac{\bar{D} - 0}{\sqrt{\frac{s_D^2}{n_D}}}$

# Two-sample test of proportions

## **The Setting:**

In two independent populations, individuals can be classified into one of two categories. You draw samples from both populations to make inference about whether the category proportions are the same in both populations

## **Example:**

“You take a sample of 100 government concentrators and 100 concentrators from other fields. Test whether being a government concentrator makes a college student more or less likely to vote.”

# Two-sample test of proportions

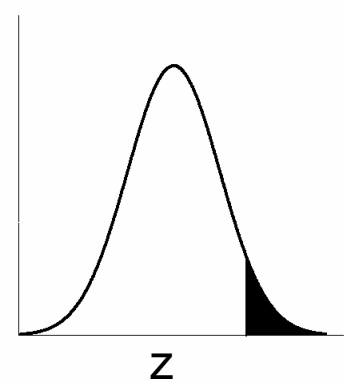
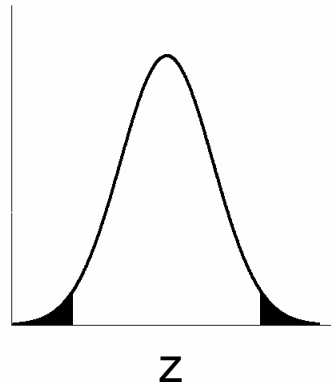
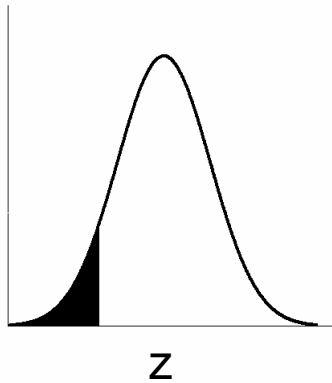
Hypothesis test:

$$H_0: p_x = p_y$$

$$H_A: p_x < p_y$$

$$H_A: p_x \neq p_y$$

$$H_A: p_x > p_y$$



The test statistic is

$$z = \frac{\hat{p}_x - \hat{p}_y}{\sqrt{\hat{p}(1 - \hat{p}) \left( \frac{1}{n_x} + \frac{1}{n_y} \right)}}$$

Where 
$$\hat{p} = \frac{n_x \hat{p}_x + n_y \hat{p}_y}{n_x + n_y}$$

# Two-sample test of proportions

Level-C Confidence Interval:

$$(p_x - p_y) = (\hat{p}_x - \hat{p}_y) \pm z * \sqrt{\frac{\hat{p}_x(1 - \hat{p}_x)}{n_x} + \frac{\hat{p}_y(1 - \hat{p}_y)}{n_y}}$$

- **Common Mistake:** Notice that in this case the standard errors of the hypothesis test and confidence interval are different!

# Chi-square test

**The Setting:** In two (or more) independent populations, individuals can be classified into two (or more) different categories. You draw samples from all populations to test whether the category an individual falls into is independent of that individual's population.

## **Example:**

“You choose a sample of 100 men and 100 women and ask subject whether s/he watches professional football. Test whether gender is independent of football viewing”



# Chi-square test

Hypothesis test (Continuing our example):

$H_0$ : gender is independent of watching football

$H_A$ : gender is not independent of watching football

- Remember that we cannot test a one-sided alternative using the chi-square test!

# Chi-square test

## **Procedure:**

1. Create the 2x2 table of observed counts
2. Create the 2x2 table of expected counts
3. Calculate the  $\chi^2$  statistic
4. Draw conclusion

# Chi-square test

- 1. Create the 2x2 table of observed counts**
- 2. Create the 2x2 table of expected counts**

$$\text{Expected count in a cell} = \frac{r \times c}{n}$$

Where  $r$  = row total for that cell,

$c$  = column total for that cell,

$n$  = total for whole table

# Chi-square test

## 3. Calculate the $\chi^2$ statistic

$$\chi^2 = \sum \frac{(\textit{Observed} - \textit{Expected})^2}{\textit{Expected}}$$

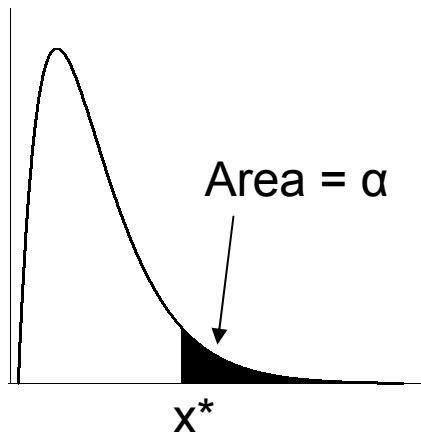
Where the summation is over all cells.

# Chi-square test

## 4. Draw conclusion

Under  $H_0$ ,  $\chi^2$  follows a chi-square distribution with  $(r-1)(c-1)$  degrees of freedom

To use the critical value method, determine the value  $x^*$  that cuts off  $\alpha$  of the tail area. Reject  $H_0$  if  $\chi^2 > x^*$ .



To use the p-value method, determine the area to the right of  $\chi^2$  on the chi-square curve. Report this p-value and draw your conclusion.

