

Inference for Distributions of Categorical Data

lesson 1

Chi-Square Goodness-of-Fit Tests

Chi-Square Goodness-of-Fit Tests

Learning Objectives

After this section, you should be able to...

- ✓ *COMPUTE* expected counts, conditional distributions, and contributions to the chi-square statistic
- ✓ *CHECK* the Random, Large sample size, and Independent conditions before performing a chi-square test
- ✓ *PERFORM* a chi-square goodness-of-fit test to determine whether sample data are consistent with a specified distribution of a categorical variable
- ✓ *EXAMINE* individual components of the chi-square statistic as part of a follow-up analysis

- Here's what the company says about the color distribution of its M&M'S Milk Chocolate Candies: On average, the M&M'S Milk Chocolate Candies will contain 13% of each of browns and reds, 14% yellows, 16% greens, 20% oranges and 24% blues.
- The **one-way table** below summarizes the data from a sample bag of M&M'S Milk Chocolate Candies. In general, one-way tables display the distribution of a categorical variable for the individuals in a sample.

Color	Blue	Orange	Green	Yellow	Red	Brown	Total
Count	9	8	12	15	10	6	60

- Since the company claims that 24% of all M&M'S Milk Chocolate Candies are blue, we could use the one-sample z test for a proportion.

$$H_0: p = 0.24$$

$$H_a: p \neq 0.24$$

where p is the true population proportion of blue M&M'S. We could then perform additional significance tests for each of the remaining colors.

× pretty inefficient ×

■ Comparing Observed and Expected Counts

We need a new kind of significance test, called
chi-square goodness-of-fit test

The null hypothesis in a chi-square goodness-of-fit test should state a claim about the distribution of a single categorical variable in the population of interest. In our example, the appropriate null hypothesis is

H_0 : The company's stated color distribution for
M&M'S Milk Chocolate Candies is correct.

The alternative hypothesis in a chi-square goodness-of-fit test is that the categorical variable does *not* have the specified distribution. In our example, the alternative hypothesis is

H_a : The company's stated color distribution for
M&M'S Milk Chocolate Candies is not correct.

■ Comparing Observed and Expected Counts

We can also write the hypotheses in symbols as:

$$H_0: p_{blue} = 0.24, p_{orange} = 0.20, p_{green} = 0.16, p_{yellow} = 0.14, p_{red} = 0.13, p_{brown} = 0.13,$$

$$H_a: \text{At least one of the } p_i\text{'s is incorrect}$$

where p_{color} is the true **population** proportion of M&M'S Milk Chocolate Candies of that color.

The idea of the chi-square goodness-of-fit test:

- Assume H_0 is true, compare the **observed counts** (from sample) with the counts that would be expected.
- The more the observed counts differ from the **expected counts**, the more evidence we have against the null hypothesis.

■ Example: Computing Expected Counts

A sample bag of M&M's milk Chocolate Candies contained 60 candies.
Calculate the expected counts for each color.

Assuming that the color distribution stated by the company is true.

Color	Observed	Expected
Blue	9	14.40
Orange	8	12.00
Green	12	9.60
Yellow	15	8.40
Red	10	7.80
Brown	6	7.80

Ho: 13% of each of browns and reds, 14% yellows, 16% greens, 20% oranges and 24% blues.

Definition:

The **chi-square statistic** is a measure of how far the observed counts are from the expected counts. The formula for the statistic is

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

where the sum is over all possible values of the categorical variable.

■ Example: Return of the M&M's

The table shows the observed and expected counts for our sample of 60 M&M's Milk Chocolate Candies. Calculate the chi-square statistic.

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

Color	Observed	Expected
Blue	9	14.40
Orange	8	12.00
Green	12	9.60
Yellow	15	8.40
Red	10	7.80
Brown	6	7.80

$$\chi^2 = \frac{(9 - 14.40)^2}{14.40} + \frac{(8 - 12.00)^2}{12.00} + \frac{(12 - 9.60)^2}{9.60} + \frac{(15 - 8.40)^2}{8.40} + \frac{(10 - 7.80)^2}{7.80} + \frac{(6 - 7.80)^2}{7.80}$$

$$\chi^2 = 2.025 + 1.333 + 0.600 + 5.186 + 0.621 + 0.415 = 10.180$$

Chi-Square Goodness-of-Fit Tests

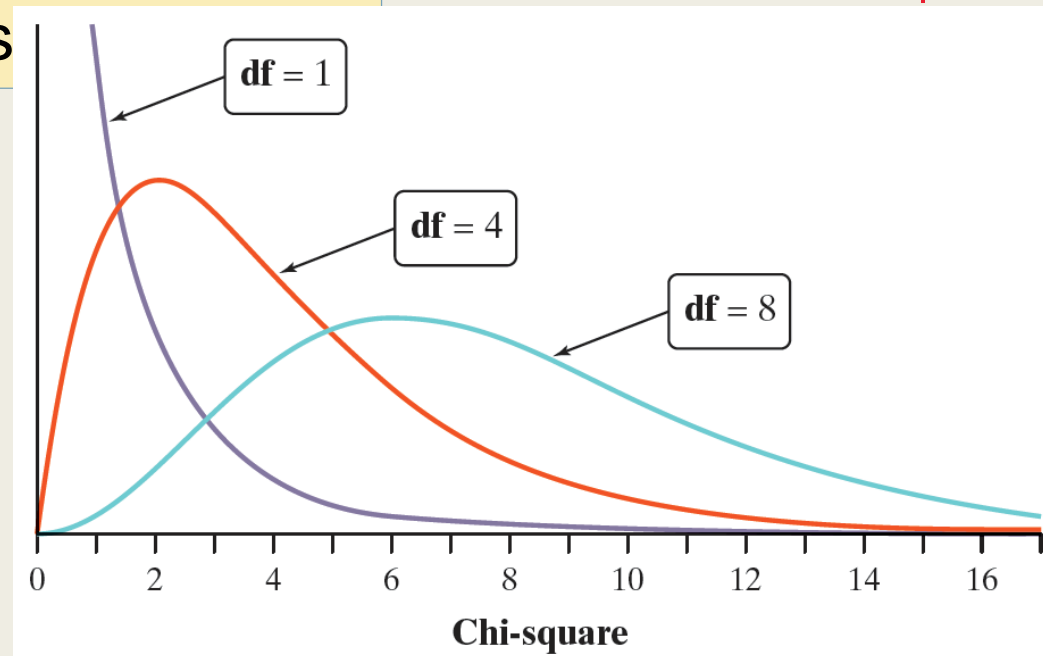
■ The Chi-Square Distributions

The Chi-Square Distributions

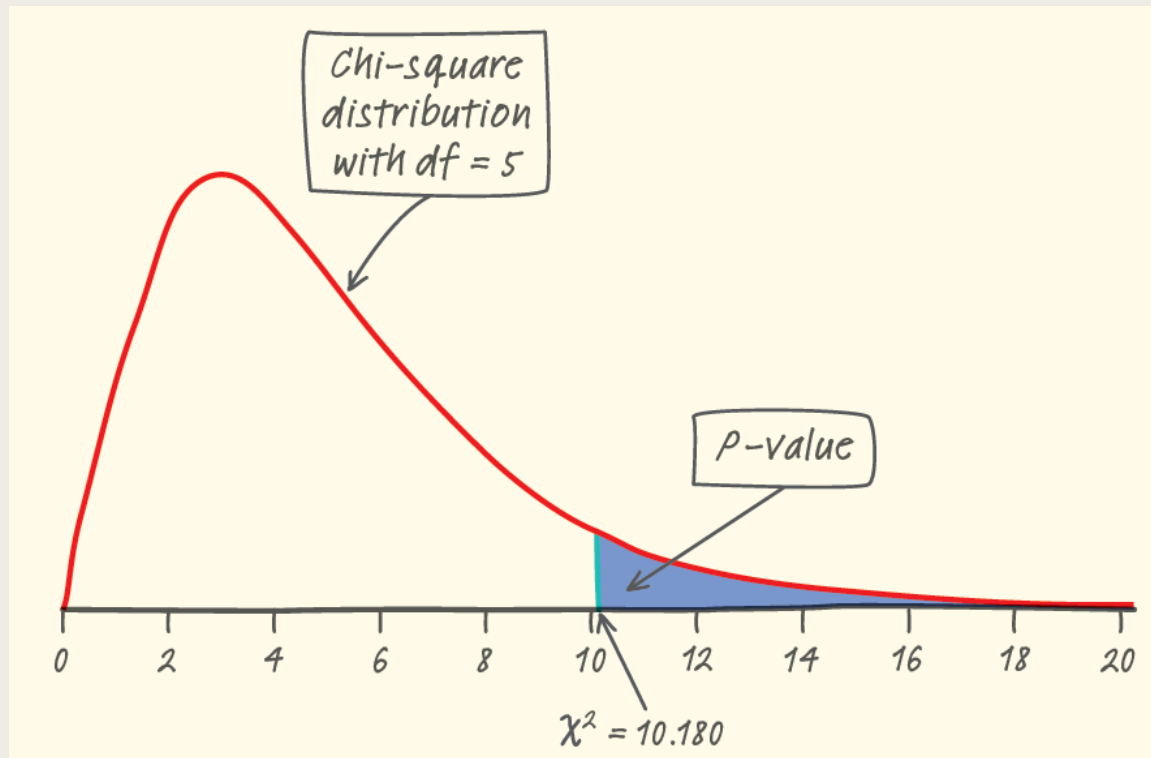
The chi-square distributions are a family of distributions that take only positive values and are skewed to the right.

The chi-square goodness-of-fit test uses the chi-square distribution with $df = n - 1$. n is the number of categories

Chi-Square Goodness-of-



■ Example: Return of the M&M's



- Since our P-value(0.07) is greater than $\alpha(0.05)$.
- Therefore, we fail to reject H_0 .
- We don't have sufficient evidence to conclude that the company's claimed color distribution is incorrect.

■ Carrying Out a Test

Conditions: Use the chi-square goodness-of-fit test when

✓ **Random**

The data come from a random sample or a randomized experiment.

✓ **Large Sample Size**

All expected counts are at least 5.

✓ **Independent**

Individual observations are independent.

When sampling without replacement, check that the population is at least 10 times as large as the sample (the 10% condition).

The Chi-Square Goodness-of-Fit Test

Suppose the k conditions are met. To determine if the distribution, expressed in terms of category, proportions, or probabilities, is a good fit to the data, we use the chi-square goodness-of-fit test.

Before we start using the chi-square goodness-of-fit test, we have two important cautions to offer.

1. The chi-square test statistic compares observed and expected *counts*. Don't try to perform calculations with the observed and expected *proportions* in each category.
2. When checking the Large Sample Size condition, be sure to examine the *expected* counts, not the observed counts.

where the sum is over all categories. The P -value is the area to the right of χ^2 under the density curve of the chi-square distribution with $k - 1$ degrees of freedom.

■ Example: When Were You Born?

Are births evenly distributed across the days of the week? The one-way table below shows the distribution of births across the days of the week in a random sample of 140 births from local records in a large city. Do these data give significant evidence that local births are not equally likely on all days of the week?

Day	Sun	Mon	Tue	Wed	Thu	Fri	Sat
Births	13	23	24	20	27	18	15

State: We want to perform a test of

H_0 : Birth days in this local area are evenly distributed across the days of the week.

H_a : Birth days in this local area are not evenly distributed across the days of the week.

The null hypothesis says that the proportions of births are the same on all days. In that case, all 7 proportions must be $1/7$. So we could also write the hypotheses as

$$H_0: p_{Sun} = p_{Mon} = p_{Tues} = \dots = p_{Sat} = 1/7.$$

H_a : At least one of the proportions is not $1/7$.

We will use $\alpha = 0.05$.

Plan: If the conditions are met, we should conduct a chi-square goodness-of-fit test.

- *Random* The data came from a random sample of local births.
- *Large Sample Size* Assuming H_0 is true, we would expect one-seventh of the births to occur on each day of the week. For the sample of 140 births, the expected count for all 7 days would be $1/7(140) = 20$ births. Since $20 \geq 5$, this condition is met.
- *Independent* Individual births in the random sample should occur independently (assuming no twins). Because we are sampling without replacement, there need to be at least $10(140) = 1400$ births in the local area. This should be the case in a large city.

■ Example: When Were You Born?

Do: Since the conditions are satisfied, we can perform a chi-square goodness-of-fit test. We begin by calculating the test statistic.

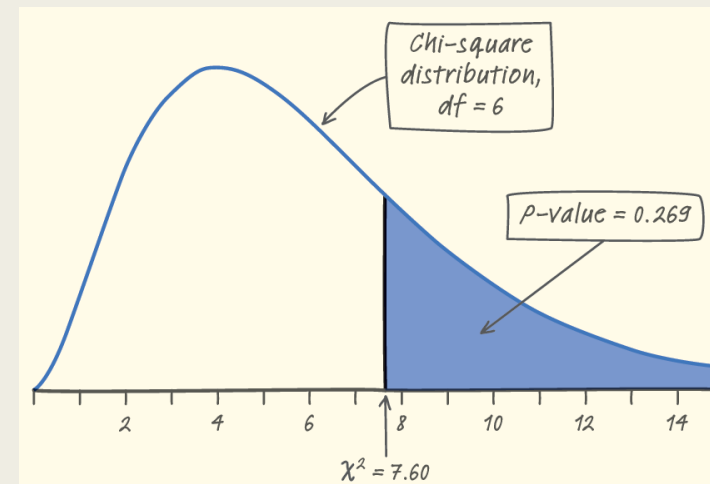
Test statistic:

$$\begin{aligned}\chi^2 &= \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} \\ &= \frac{(13-20)^2}{20} + \frac{(23-20)^2}{20} + \frac{(24-20)^2}{20} + \frac{(20-20)^2}{20} \\ &\quad + \frac{(27-20)^2}{20} + \frac{(18-20)^2}{20} + \frac{(15-20)^2}{20} \\ &= 2.45 + 0.45 + 0.80 + 0.00 + 2.45 + 0.20 + 1.25 \\ &= 7.60\end{aligned}$$

P-Value:

Using Table C: $\chi^2 = 7.60$ is less than the smallest entry in the $df = 6$ row, which corresponds to tail area 0.25. The P -value is therefore greater than 0.25.

Using technology: We can find the exact P -value with a calculator: $\chi^2\text{cdf}(7.60, 1000, 6) = 0.269$.



Conclude: Because the P -value, 0.269, is greater than $\alpha = 0.05$, we fail to reject H_0 . These 140 births don't provide enough evidence to say that all local births in this area are not evenly distributed across the days of the week.