



Test of Independence

By:
OpenStaxCollege

Tests of independence involve using a contingency table of observed (data) values.

The test statistic for a test of independence is similar to that of a goodness-of-fit test:

$$\sum_{(i \cdot j)} \frac{(O - E)^2}{E}$$

where:

- O = observed values
- E = expected values
- i = the number of rows in the table
- j = the number of columns in the table

There are $i \cdot j$ terms of the form $\frac{(O - E)^2}{E}$.

A test of independence determines whether two factors are independent or not. You first encountered the term independence in [Probability Topics](#). As a review, consider the following example.

Note

The expected value for each cell needs to be at least five in order for you to use this test.

Suppose A = a speeding violation in the last year and B = a cell phone user while driving. If A and B are independent then $P(A \text{ AND } B) = P(A)P(B)$. $A \text{ AND } B$ is the event that a driver received a speeding violation last year and also used a cell phone while driving. Suppose, in a study of drivers who received speeding violations in the last year, and who used cell phone while driving, that 755 people were surveyed. Out of the 755, 70 had a speeding violation and 685 did not; 305 used cell phones while driving and 450 did not.

Let y = expected number of drivers who used a cell phone while driving and received speeding violations.

Test of Independence

If A and B are independent, then $P(A \text{ AND } B) = P(A)P(B)$. By substitution,

$$\frac{y}{755} = \left(\frac{70}{755}\right)\left(\frac{305}{755}\right)$$

$$\text{Solve for } y: y = \frac{(70)(305)}{755} = 28.3$$

About 28 people from the sample are expected to use cell phones while driving and to receive speeding violations.

In a test of independence, we state the null and alternative hypotheses in words. Since the contingency table consists of **two factors**, the null hypothesis states that the factors are **independent** and the alternative hypothesis states that they are **not independent (dependent)**. If we do a test of independence using the example, then the null hypothesis is:

H_0 : Being a cell phone user while driving and receiving a speeding violation are independent events.

If the null hypothesis were true, we would expect about 28 people to use cell phones while driving and to receive a speeding violation.

The test of independence is always right-tailed because of the calculation of the test statistic. If the expected and observed values are not close together, then the test statistic is very large and way out in the right tail of the chi-square curve, as it is in a goodness-of-fit.

The number of degrees of freedom for the test of independence is:

$$df = (\text{number of columns} - 1)(\text{number of rows} - 1)$$

The following formula calculates the **expected number** (E):

$$E = \frac{(\text{row total})(\text{column total})}{\text{total number surveyed}}$$

Try It

A sample of 300 students is taken. Of the students surveyed, 50 were music students, while 250 were not. Ninety-seven were on the honor roll, while 203 were not. If we assume being a music student and being on the honor roll are independent events, what is the expected number of music students who are also on the honor roll?

About 16 students are expected to be music students and on the honor roll.

Test of Independence

In a volunteer group, adults 21 and older volunteer from one to nine hours each week to spend time with a disabled senior citizen. The program recruits among community college students, four-year college students, and nonstudents. In [\[link\]](#) is a **sample** of the adult volunteers and the number of hours they volunteer per week.

Number of Hours Worked Per Week by Volunteer Type (Observed)The table contains **observed (O)** values (data).

Type of Volunteer	1–3 Hours	4–6 Hours	7–9 Hours	Row Total
Community College Students	111	96	48	255
Four-Year College Students	96	133	61	290
Nonstudents	91	150	53	294
Column Total	298	379	162	839

Is the number of hours volunteered **independent** of the type of volunteer?

The **observed table** and the question at the end of the problem, "Is the number of hours volunteered independent of the type of volunteer?" tell you this is a test of independence. The two factors are **number of hours volunteered** and **type of volunteer**. This test is always right-tailed.

H_0 : The number of hours volunteered is **independent** of the type of volunteer.

H_a : The number of hours volunteered is **dependent** on the type of volunteer.

The expected result are in [\[link\]](#).

Number of Hours Worked Per Week by Volunteer Type
(Expected)The table contains **expected (E)** values (data).

Type of Volunteer	1-3 Hours	4-6 Hours	7-9 Hours
Community College Students	90.57	115.19	49.24
Four-Year College Students	103.00	131.00	56.00
Nonstudents	104.42	132.81	56.77

For example, the calculation for the expected frequency for the top left cell is

$$E = \frac{(\text{row total})(\text{column total})}{\text{total number surveyed}} = \frac{(255)(298)}{839} = 90.57$$

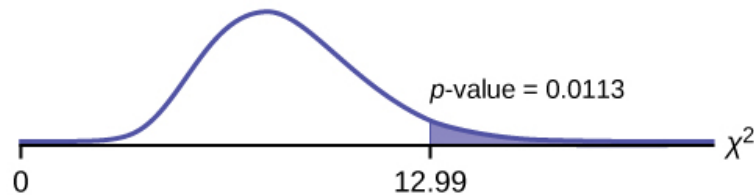
Test of Independence

Calculate the test statistic: $\chi^2 = 12.99$ (calculator or computer)

Distribution for the test: χ^2_4

$$df = (3 \text{ columns} - 1)(3 \text{ rows} - 1) = (2)(2) = 4$$

Graph:



Probability statement: $p\text{-value} = P(\chi^2 > 12.99) = 0.0113$

Compare α and the p -value: Since no α is given, assume $\alpha = 0.05$. $p\text{-value} = 0.0113$. $\alpha > p\text{-value}$.

Make a decision: Since $\alpha > p\text{-value}$, reject H_0 . This means that the factors are not independent.

Conclusion: At a 5% level of significance, from the data, there is sufficient evidence to conclude that the number of hours volunteered and the type of volunteer are dependent on one another.

For the example in [\[link\]](#), if there had been another type of volunteer, teenagers, what would the degrees of freedom be?

Press the MATRX key and arrow over to EDIT. Press 1: [A]. Press 3 ENTER 3 ENTER. Enter the table values by row from [\[link\]](#). Press ENTER after each. Press 2nd QUIT. Press STAT and arrow over to TESTS. Arrow down to C: χ^2 -TEST. Press ENTER. You should see Observed: [A] and Expected: [B]. Arrow down to Calculate. Press ENTER. The test statistic is 12.9909 and the $p\text{-value} = 0.0113$. Do the procedure a second time, but arrow down to Draw instead of calculate.

Try It

The Bureau of Labor Statistics gathers data about employment in the United States. A sample is taken to calculate the number of U.S. citizens working in one of several industry sectors over time. [\[link\]](#) shows the results:

Test of Independence

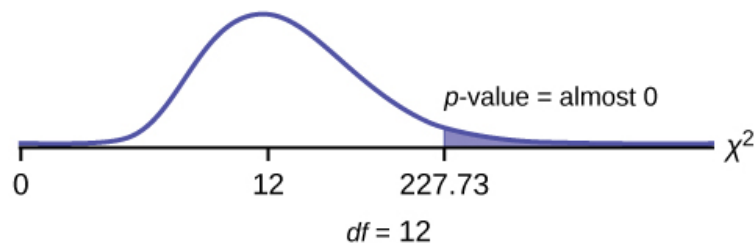
Industry Sector	2000	2010	2020	Total
Nonagriculture wage and salary	13,243	13,044	15,018	41,305
Goods-producing, excluding agriculture	2,457	1,771	1,950	6,178
Services-providing	10,786	11,273	13,068	35,127
Agriculture, forestry, fishing, and hunting	240	214	201	655
Nonagriculture self-employed and unpaid family worker	931	894	972	2,797
Secondary wage and salary jobs in agriculture and private household industries	14	11	11	36
Secondary jobs as a self-employed or unpaid family worker	196	144	152	492
Total	27,867	27,351	31,372	86,590

We want to know if the change in the number of jobs is independent of the change in years. State the null and alternative hypotheses and the degrees of freedom.

H_0 : The number of jobs is independent of the year.

H_a : The number of jobs is dependent on the year.

$df = 12$



Press the **MATRX** key and arrow over to **EDIT**. Press **1** : [A]. Press **3** **ENTER** **3** **ENTER**. Enter the table values by row. Press **ENTER** after each. Press **2nd** **QUIT**. Press **STAT** and arrow over to **TESTS**. Arrow down to **C**: χ^2 -TEST. Press **ENTER**. You should see **Observed**: [A] and **Expected**: [B]. Arrow down to **Calculate**. Press **ENTER**. The test statistic is 227.73 and the p -value = $5.90E - 42 = 0$. Do the procedure a second time but arrow down to **Draw** instead of **calculate**.

De Anza College is interested in the relationship between anxiety level and the need to succeed in school. A random sample of 400 students took a test that measured anxiety

Test of Independence

level and need to succeed in school. [\[link\]](#) shows the results. De Anza College wants to know if anxiety level and need to succeed in school are independent events.

Need to Succeed in School vs. Anxiety Level

Need to Succeed in School	High Anxiety	Med-high Anxiety	Medium Anxiety	Med-low Anxiety	Low Anxiety	Row Total
High Need	35	42	53	15	10	155
Medium Need	18	48	63	33	31	193
Low Need	4	5	11	15	17	52
Column Total	57	95	127	63	58	400

a. How many high anxiety level students are expected to have a high need to succeed in school?

a. The column total for a high anxiety level is 57. The row total for high need to succeed in school is 155. The sample size or total surveyed is 400.

$$E = \frac{(\text{row total})(\text{column total})}{\text{total surveyed}} = \frac{155 \cdot 57}{400} = 22.09$$

The expected number of students who have a high anxiety level and a high need to succeed in school is about 22.

b. If the two variables are independent, how many students do you expect to have a low need to succeed in school and a med-low level of anxiety?

b. The column total for a med-low anxiety level is 63. The row total for a low need to succeed in school is 52. The sample size or total surveyed is 400.

c. $E = \frac{(\text{row total})(\text{column total})}{\text{total surveyed}} = \underline{\hspace{2cm}}$

c. $E = \frac{(\text{row total})(\text{column total})}{\text{total surveyed}} = 8.19$

d. The expected number of students who have a med-low anxiety level and a low need to succeed in school is about .

d. 8

Try It

Test of Independence

Refer back to the information in [\[link\]](#). How many service providing jobs are there expected to be in 2020? How many nonagriculture wage and salary jobs are there expected to be in 2020?

12,727, 14,965

References

DiCamilo, Mark, Mervin Field, “Most Californians See a Direct Linkage between Obesity and Sugary Sodas. Two in Three Voters Support Taxing Sugar-Sweetened Beverages If Proceeds are Tied to Improving School Nutrition and Physical Activity Programs.” The Field Poll, released Feb. 14, 2013. Available online at <http://field.com/fieldpollonline/subscribers/RIs2436.pdf> (accessed May 24, 2013).

Harris Interactive, “Favorite Flavor of Ice Cream.” Available online at <http://www.statisticbrain.com/favorite-flavor-of-ice-cream> (accessed May 24, 2013)

“Youngest Online Entrepreneurs List.” Available online at <http://www.statisticbrain.com/youngest-online-entrepreneur-list> (accessed May 24, 2013).

Chapter Review

To assess whether two factors are independent or not, you can apply the test of independence that uses the chi-square distribution. The null hypothesis for this test states that the two factors are independent. The test compares observed values to expected values. The test is right-tailed. Each observation or cell category must have an expected value of at least 5.

Formula Review

Test of Independence

- The number of degrees of freedom is equal to (number of columns - 1)(number of rows - 1).
- The test statistic is $\sum_{(i \cdot j)} \frac{(O - E)^2}{E}$ where O = observed values, E = expected values, i = the number of rows in the table, and j = the number of columns in the table.
- If the null hypothesis is true, the expected number $E = \frac{(\text{row total})(\text{column total})}{\text{total surveyed}}$.

Determine the appropriate test to be used in the next three exercises.

Test of Independence

A pharmaceutical company is interested in the relationship between age and presentation of symptoms for a common viral infection. A random sample is taken of 500 people with the infection across different age groups.

a test of independence

The owner of a baseball team is interested in the relationship between player salaries and team winning percentage. He takes a random sample of 100 players from different organizations.

A marathon runner is interested in the relationship between the brand of shoes runners wear and their run times. She takes a random sample of 50 runners and records their run times as well as the brand of shoes they were wearing.

a test of independence

Use the following information to answer the next seven exercises: Transit Railroads is interested in the relationship between travel distance and the ticket class purchased. A random sample of 200 passengers is taken. [\[link\]](#) shows the results. The railroad wants to know if a passenger's choice in ticket class is independent of the distance they must travel.

Traveling Distance	Third class	Second class	First class	Total
1–100 miles	21	14	6	41
101–200 miles	18	16	8	42
201–300 miles	16	17	15	48
301–400 miles	12	14	21	47
401–500 miles	6	6	10	22
Total	73	67	60	200

State the hypotheses.

H_0 : _____

H_a : _____

df = _____

Test of Independence

How many passengers are expected to travel between 201 and 300 miles and purchase second-class tickets?

How many passengers are expected to travel between 401 and 500 miles and purchase first-class tickets?

6.6

What is the test statistic?

What is the p -value?

0.0435

What can you conclude at the 5% level of significance?

Use the following information to answer the next eight exercises: An article in the New England Journal of Medicine, discussed a study on smokers in California and Hawaii. In one part of the report, the self-reported ethnicity and smoking levels per day were given. Of the people smoking at most ten cigarettes per day, there were 9,886 African Americans, 2,745 Native Hawaiians, 12,831 Latinos, 8,378 Japanese Americans and 7,650 whites. Of the people smoking 11 to 20 cigarettes per day, there were 6,514 African Americans, 3,062 Native Hawaiians, 4,932 Latinos, 10,680 Japanese Americans, and 9,877 whites. Of the people smoking 21 to 30 cigarettes per day, there were 1,671 African Americans, 1,419 Native Hawaiians, 1,406 Latinos, 4,715 Japanese Americans, and 6,062 whites. Of the people smoking at least 31 cigarettes per day, there were 759 African Americans, 788 Native Hawaiians, 800 Latinos, 2,305 Japanese Americans, and 3,970 whites.

Complete the table.

Smoking Levels by Ethnicity (Observed)

Smoking Level Per Day	African American	Native Hawaiian	Latino	Japanese Americans	White	TOTALS
1-10						
11-20						
21-30						
31+						
TOTALS						

Test of Independence

Smoking Level Per Day	African American	Native Hawaiian	Latino	Japanese Americans	White	Totals
1-10	9,886	2,745	12,831	8,378	7,650	41,490
11-20	6,514	3,062	4,932	10,680	9,877	35,065
21-30	1,671	1,419	1,406	4,715	6,062	15,273
31+	759	788	800	2,305	3,970	8,622
Totals	18,830	8,014	19,969	26,078	27,559	10,0450

State the hypotheses.

H_0 : _____

H_a : _____

Enter expected values in [\[link\]](#). Round to two decimal places.

Calculate the following values:

Smoking Level Per Day	African American	Native Hawaiian	Latino	Japanese Americans	White
1-10	7777.57	3310.11	8248.02	10771.29	11383.01
11-20	6573.16	2797.52	6970.76	9103.29	9620.27
21-30	2863.02	1218.49	3036.20	3965.05	4190.23
31+	1616.25	687.87	1714.01	2238.37	2365.49

df = _____

χ^2 test statistic = _____

10,301.8

p -value = _____

Is this a right-tailed, left-tailed, or two-tailed test? Explain why.

right

Graph the situation. Label and scale the horizontal axis. Mark the mean and test statistic. Shade in the region corresponding to the p -value.

Test of Independence



State the decision and conclusion (in a complete sentence) for the following preconceived levels of α .

$$\alpha = 0.05$$

1. Decision: _____
 2. Reason for the decision: _____
 3. Conclusion (write out in a complete sentence): _____
1. Reject the null hypothesis.
 2. p -value $<$ alpha
 3. There is sufficient evidence to conclude that smoking level is dependent on ethnic group.

$$\alpha = 0.01$$

1. Decision: _____
2. Reason for the decision: _____
3. Conclusion (write out in a complete sentence): _____

Homework

For each problem, use a solution sheet to solve the hypothesis test problem. Go to [Appendix E](#) for the chi-square solution sheet. Round expected frequency to two decimal places.

A recent debate about where in the United States skiers believe the skiing is best prompted the following survey. Test to see if the best ski area is independent of the level of the skier.

U.S. Ski Area	Beginner	Intermediate	Advanced
Tahoe	20	30	40

Test of Independence

U.S. Ski Area	Beginner	Intermediate	Advanced
Utah	10	30	60
Colorado	10	40	50

Car manufacturers are interested in whether there is a relationship between the size of car an individual drives and the number of people in the driver's family (that is, whether car size and family size are independent). To test this, suppose that 800 car owners were randomly surveyed with the results in [\[link\]](#). Conduct a test of independence.

Family Size	Sub & Compact	Mid-size	Full-size	Van & Truck
1	20	35	40	35
2	20	50	70	80
3–4	20	50	100	90
5+	20	30	70	70

1. H_0 : Car size is independent of family size.
2. H_a : Car size is dependent on family size.
3. $df = 9$
4. chi-square distribution with $df = 9$
5. test statistic = 15.8284
6. p -value = 0.0706
7. Check student's solution.
8.
 1. Alpha: 0.05
 2. Decision: Do not reject the null hypothesis.
 3. Reason for decision: p -value > alpha
 4. Conclusion: At the 5% significance level, there is insufficient evidence to conclude that car size and family size are dependent.

College students may be interested in whether or not their majors have any effect on starting salaries after graduation. Suppose that 300 recent graduates were surveyed as to their majors in college and their starting salaries after graduation. [\[link\]](#) shows the data. Conduct a test of independence.

Major	< \$50,000	\$50,000 – \$68,999	\$69,000 +
English	5	20	5
Engineering	10	30	60
Nursing	10	15	15

Test of Independence

Major	< \$50,000	\$50,000 – \$68,999	\$69,000 +
Business	10	20	30
Psychology	20	30	20

Some travel agents claim that honeymoon hot spots vary according to age of the bride. Suppose that 280 recent brides were interviewed as to where they spent their honeymoons. The information is given in [\[link\]](#). Conduct a test of independence.

Location	20–29	30–39	40–49	50 and over
Niagara Falls	15	25	25	20
Poconos	15	25	25	10
Europe	10	25	15	5
Virgin Islands	20	25	15	5

1. H_0 : Honeymoon locations are independent of bride's age.
2. H_a : Honeymoon locations are dependent on bride's age.
3. $df = 9$
4. chi-square distribution with $df = 9$
5. test statistic = 15.7027
6. p -value = 0.0734
7. Check student's solution.
8.
 1. Alpha: 0.05
 2. Decision: Do not reject the null hypothesis.
 3. Reason for decision: p -value > alpha
 4. Conclusion: At the 5% significance level, there is insufficient evidence to conclude that honeymoon location and bride age are dependent.

A manager of a sports club keeps information concerning the main sport in which members participate and their ages. To test whether there is a relationship between the age of a member and his or her choice of sport, 643 members of the sports club are randomly selected. Conduct a test of independence.

Sport	18 - 25	26 - 30	31 - 40	41 and over
racquetball	42	58	30	46
tennis	58	76	38	65
swimming	72	60	65	33

Test of Independence

A major food manufacturer is concerned that the sales for its skinny french fries have been decreasing. As a part of a feasibility study, the company conducts research into the types of fries sold across the country to determine if the type of fries sold is independent of the area of the country. The results of the study are shown in [\[link\]](#). Conduct a test of independence.

Type of Fries	Northeast	South	Central	West
skinny fries	70	50	20	25
curly fries	100	60	15	30
steak fries	20	40	10	10

1. H_0 : The types of fries sold are independent of the location.
2. H_a : The types of fries sold are dependent on the location.
3. $df = 6$
4. chi-square distribution with $df = 6$
5. test statistic = 18.8369
6. p -value = 0.0044
7. Check student's solution.
8.
 1. Alpha: 0.05
 2. Decision: Reject the null hypothesis.
 3. Reason for decision: p -value < alpha
 4. Conclusion: At the 5% significance level, There is sufficient evidence that types of fries and location are dependent.

According to Dan Lenard, an independent insurance agent in the Buffalo, N.Y. area, the following is a breakdown of the amount of life insurance purchased by males in the following age groups. He is interested in whether the age of the male and the amount of life insurance purchased are independent events. Conduct a test for independence.

Age of Males	None	< \$200,000	\$200,000–\$400,000	\$401,001–\$1,000,000	\$1,000,001+
20–29	40	15	40	0	5
30–39	35	5	20	20	10
40–49	20	0	30	0	30
50+	40	30	15	15	10

Test of Independence

Suppose that 600 thirty-year-olds were surveyed to determine whether or not there is a relationship between the level of education an individual has and salary. Conduct a test of independence.

Annual Salary	Not a high school graduate	High school graduate	College graduate	Masters or doctorate
< \$30,000	15	25	10	5
\$30,000–\$40,000	20	40	70	30
\$40,000–\$50,000	10	20	40	55
\$50,000–\$60,000	5	10	20	60
\$60,000+	0	5	10	150

1. H_0 : Salary is independent of level of education.
2. H_a : Salary is dependent on level of education.
3. $df = 12$
4. chi-square distribution with $df = 12$
5. test statistic = 255.7704
6. p -value = 0
7. Check student's solution.

8. Alpha: 0.05

Decision: Reject the null hypothesis.

Reason for decision: p -value < alpha

Conclusion: At the 5% significance level, there is sufficient evidence to conclude that salary and level of education are dependent.

Read the statement and decide whether it is true or false.

The number of degrees of freedom for a test of independence is equal to the sample size minus one.

The test for independence uses tables of observed and expected data values.

true

The test to use when determining if the college or university a student chooses to attend is related to his or her socioeconomic status is a test for independence.

Test of Independence

In a test of independence, the expected number is equal to the row total multiplied by the column total divided by the total surveyed.

true

An ice cream maker performs a nationwide survey about favorite flavors of ice cream in different geographic areas of the U.S. Based on [\[link\]](#), do the numbers suggest that geographic location is independent of favorite ice cream flavors? Test at the 5% significance level.

U.S. region/ Flavor	Strawberry	Chocolate	Vanilla	Rocky Road	Mint Chocolate Chip	Pistachio	Row total
West	12	21	22	19	15	8	97
Midwest	10	32	22	11	15	6	96
East	8	31	27	8	15	7	96
South	15	28	30	8	15	6	102
Column Total	45	112	101	46	60	27	391

[\[link\]](#) provides a recent survey of the youngest online entrepreneurs whose net worth is estimated at one million dollars or more. Their ages range from 17 to 30. Each cell in the table illustrates the number of entrepreneurs who correspond to the specific age group and their net worth. Are the ages and net worth independent? Perform a test of independence at the 5% significance level.

Age Group\ Net Worth Value (in millions of US dollars)	1-5	6-24	≥ 25	Row Total
17-25	8	7	5	20
26-30	6	5	9	20
Column Total	14	12	14	40

1. H_0 : Age is independent of the youngest online entrepreneurs' net worth.
2. H_a : Age is dependent on the net worth of the youngest online entrepreneurs.
3. $df = 2$
4. chi-square distribution with $df = 2$
5. test statistic = 1.76
6. p -value 0.4144

Test of Independence

7. Check student's solution.
8.
 1. Alpha: 0.05
 2. Decision: Do not reject the null hypothesis.
 3. Reason for decision: $p\text{-value} > \alpha$
 4. Conclusion: At the 5% significance level, there is insufficient evidence to conclude that age and net worth for the youngest online entrepreneurs are dependent.

A 2013 poll in California surveyed people about taxing sugar-sweetened beverages. The results are presented in [\[link\]](#), and are classified by ethnic group and response type. Are the poll responses independent of the participants' ethnic group? Conduct a test of independence at the 5% significance level.

Opinion/ Ethnicity	Asian- American	White/Non- Hispanic	African- American	Latino	Row Total
Against tax	48	433	41	160	628
In Favor of tax	54	234	24	147	459
No opinion	16	43	16	19	84
Column Total	118	710	71	272	1171