

Categorical Data Analysis

[See online here](#)

Categorical data refers to that type of data which can be classified into groups. There are various tools and techniques that are used to analyze categorical data: pie-charts, bar charts and two-way tables (for stand-alone categorical data analysis); Chi-square test and Diagnostic-odds-ratio test (for two categorical variable relationship analysis); logistic, probit and OLS (for categorical dependent variable analysis); and simple regression (for categorical independent variable analysis).

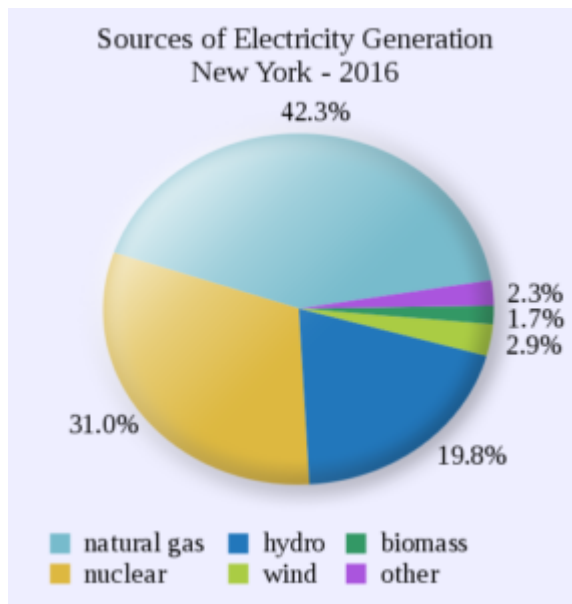


Introduction

Categorical data refers to that type of data that can be **classified into groups**. It is also called nominal data, qualitative data, etc. **Gender, marital status, income group, etc.** are some of the examples of categorical data. Consider the figure below which shows a categorical variable 'source of electricity generation. The data could be classified into six groups:

1. Natural gas
2. Hydro
3. Biomass
4. Nuclear
5. Wind
6. Other

Hence, the source of electricity generation is called categorical data.



[Image:](#) "Pie chart showing sources of electricity generated in New York" by Aflafa1. License: [CC0 1.0](#)

There are various tools and techniques that are used to analyze categorical data.

1. First of all, there are tools and techniques which are used to analyze categorical variables **on a stand-alone basis**.
2. Then there are tools and techniques which are used **when only the dependent variable is a categorical variable**.
3. Also, there are tools and techniques which are used when **only the independent variable is a categorical variable**.
4. Finally, there are tools and techniques which are used **when both the independent variables, as well as the dependent variables, are categorical variables**.

Tools and Techniques Used to Analyze Categorical Data on a Standalone Basis

The tools that are used to analyze categorical data include pie charts, bar charts, and two-by-two tables. Consider the following categorical data 'population by income group':

- Low income: 50.00
- Middle income: 130.00
- High income: 30.00

Furthermore, consider the following data 'Sales of different fruits':

- Grapes: 385000.00
- Apples: 874585.00
- Bananas: 45575.00

This categorical data could be analyzed using a graphical toolbar chart as follows:

$$\sum \text{all cells } (\text{observed count} - \text{expected count})^2 / (\text{expected count})$$

This bar chart helped us in determining the fruit which has the highest sales and the fruit

which has the lowest sales. Apples have the highest sales as it has the highest bar and bananas have the lowest sales as it has the lowest bar; thus, the fruit seller must procure more apples and fewer bananas.

Moreover, consider the following categorical data:

There are a total of 619 people in a class:

- 481 of them are male
- 138 of them are female
- 93 of them hire a home tutor
- 526 of them do not hire a home tutor
- 77 of them are male and hire a home tutor
- 16 of them are female and hire a home tutor
- 404 of them are male and do not hire a home tutor
- 122 of them are female and do not hire a home tutor

This categorical data could be depicted in a meaningful manner using a two-by-two table as:

	Hired home tutor?	Hired home tutor?	
	Yes	No	Total
Male	77	404	481
Female	16	122	138
Total	93	526	619

Tools and Techniques Used When Both Dependent and Independent Variable Are Categorical

Chi-square test for independence

One of the most common techniques used for analyzing the relationship between two categorical variables is the Chi-square test for independence. Chi is a Greek letter that looks like this: χ , so the test is sometimes referred to as The χ^2 test for independence.

The Chi-square test of independence would be explained using an example. Consider the following two-by-two table:

	Hired home tutor?	Hired home tutor?	
	Yes	No	Total
Male	77	404	481
Female	16	122	138
Total	93	526	619

Suppose we want to analyze the relationship between gender and the decision of hiring a home tutor. We would conduct the Chi-square test of independence as follows in order to answer this question.

Step 1: Defining the hypothesis

H0: Null Hypothesis: There is no relationship between the two categorical variables 'gender' and 'hired a home tutor,' i.e., they are independent.

HA: Alternative Hypothesis: There is a relationship between the two categorical variables

'gender' and 'hired a home tutor,' i.e., they are not independent.

Step 2: Compute the Chi-square (χ^2) test statistic

$$\chi^2 = \sum \text{all cells } (\text{observed count} - \text{expected count})^2 / (\text{expected count})$$

$$\text{Expected count} = (\text{Row total} * \text{Column total}) / \text{Grand total}$$

Using the above formula, the expected counts for each cell are calculated as follows:

Expected counts

	Hired home tutor?		
	Yes	No	Total
Male	$(93 * 481) / 619 = 72.3$	$(526 * 481) / 619 = 408.7$	481
Female	$(93 * 138) / 619 = 20.7$	$(526 * 138) / 619 = 117.3$	138
Total	93	526	619

We already know the **observed counts**:

	Hired home tutor?		
	Yes	No	Total
Male	77	404	481
Female	16	122	138
Total	93	526	619

Thus, the χ^2 test statistic is:

$$((77 - 72.3)^2 / 72.3) + ((404 - 408.7)^2 / 408.7) + ((16 - 20.7)^2 / 20.7) + ((122 - 117.3)^2 / 117.3) = 0.306 + 0.054 + 1.067 + 0.188 = 1.62.$$

Step 3: Finding the p-value

The p-value for the chi-square test for independence is the **probability of getting counts like those observed**, assuming that the two variables are not related (which is what is claimed by the null hypothesis). The smaller the p-value, the more surprising it would be to get counts like we did if the null hypothesis were true.

Technically, the p-value is the **probability of observing χ^2 at least as large as the one observed**. Using statistical software, we find that the p-value for this test is 0.201.

Chi-Square Test: Yes, No

	Yes	No	Total
Male	77 72.27 0.310	404 408.73 0.055	481
Female	16 20.73 1.081	122 117.27 0.191	138
Total	93	526	619

- a. Expected counts printed below are observed, counts.
- b. Chi-Square contributions printed below are expected, counts.

Chi-Square = 1.637, DF = 1, P-Value = 0.201

The p-value is higher than 0.05 and, thus, we fail to reject the null hypothesis at the 5% significance level; that is, **gender and hiring a home tutor are independent**

variables.

Diagnostic odds ratio.

It is mostly **used to test the effectiveness of a medicinal disease diagnostic test.**

The Diagnostic odds ratio test would be explained using an example. Consider the following two-by-two table:

		Actual Condition	Actual Condition
		Positive	Negative
Test Outcome	Positive	44	23
Test Outcome	Negative	6	96

That is, there were 44 patients who actually had the disease and the diagnostic test correctly revealed that they have the disease. There were 6 patients who actually had the disease but the diagnostic test incorrectly revealed that they do not have the disease. There were 23 patients who actually did not have the disease but the diagnostic test incorrectly revealed that they had the disease. Finally, there were 96 patients who did not have the disease and the diagnostic test correctly revealed that they do not have the disease.

We want to analyze the effectiveness of this diagnostic test. In order to do so, we would conduct the Diagnostic odds ratio test as follows:

Step 1: Defining the hypothesis

H₀: Null Hypothesis: The diagnostic test is not effective/accurate.

H_A: Alternative Hypothesis: The diagnostic test is effective/accurate.

Step 2: Compute the Diagnostic odds ratio

D. O. R = (True positives / False positive) / (False negatives/True negatives) = (44 / 23) / (6 / 96) = 30.6

Step 3: Making the conclusion

The D. O. R of 30.6 is greater than 1 and thus we can safely conclude that the diagnostic test is effective/accurate.

Tools and Techniques Used When Only the Dependent Variable Is Categorical

OLS regression

OLS regression **reveals the probability of selecting an option.** This would be explained using an example.

Suppose that the dependent variable is 'loan default'; that is, a categorical variable which could take only 2 values 'Yes,' if the person defaults, and 'No,' if the person does not default. The independent variables would be the income level of the person, job security of the person, etc.

Now, suppose we ran an OLS regression on the data and found out that the coefficient of independent variable 'income level' is 0.45. This would imply that if the income level rises by 1 unit, then the probability of the person defaulting on his loan would rise by 0.45.

Logistic regression

The logistic regression **reveals the log of the odds ratio of selecting an option**. This would be explained using the same example that was used for explaining OLS regression.

Now, suppose we ran a logistic regression on the data and found out that the coefficient of the independent variable 'income level' is 0.30. This would imply that if the income level rises by 1 unit, then the log of the odds ratio of the person defaulting on his loan would rise by 0.30.

The logistic distribution is shown as follows:

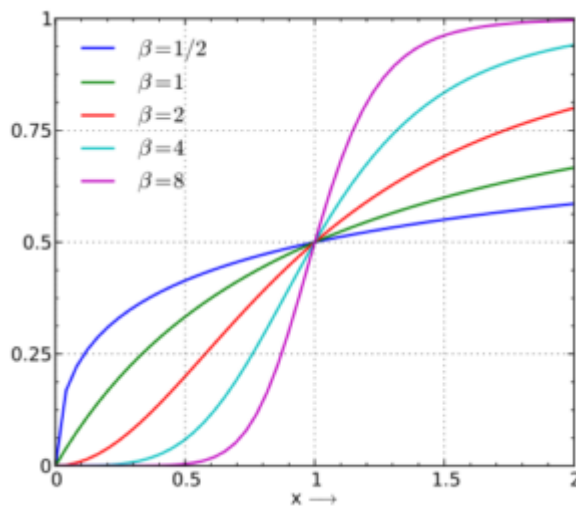


Image: "log-logistic distribution function." by Qwfp. License: [CC BY-SA 3.0](#)

The assumptions of the logistic model are as follows.

The logistic model accepts that **information is case-particular**; that is, every free factor has a solitary incentive for each case. The logistic model likewise expects that the **response variable can't be impeccably anticipated from the autonomous factors for any case**.

Similarly, as with different sorts of relapse, there is no requirement for the autonomous factors to be measurably free from each other (dissimilar to, for instance, in a gullible Bayes classifier). Be that as it may, **collinearity is thought to be generally low**, as it ends up plainly hard to separate between the effect of a few factors if this isn't the case.

In the event that the logistic is utilized to demonstrate decisions, it depends on the supposition of freedom of unimportant options, which isn't generally alluring. This suspicion expresses that the **chances of leaning toward one class over another don't rely upon the nearness or non-appearance of other "insignificant" choices**.

For instance, the relative probabilities of taking an auto or transport to work don't change if a bike is included as an extra plausibility. This enables the decision of K contrasting options to be demonstrated as an arrangement of K-1 free independent decisions, in which one option is picked as a "turn" and the other K-1 thought about against it, each one in turn.

The IIA speculation is **center speculation in the normal decision hypothesis**;

however, various investigations in brain science demonstrates that people frequently abuse this supposition when settling on decisions.

A case of an issue-case emerges if decisions incorporate an auto and a blue transport. Assume the chances proportion between the two is 1:1. Presently, if the choice of red transport is presented, a man might be aloof between a red and blue transport, and subsequently may show an auto : blue transport : red transport chances proportion of 1 : 0.5 : 0.5, accordingly keeping up a 1 : 1 proportion of auto : any transport while receiving a changed auto : blue transport proportion of 1 : 0.5.

Here the red transport choice was not in reality insignificant, on the grounds that red transport was an ideal substitute for blue transport. On the off chance that the logistic is utilized to demonstrate decisions, it might, in a few circumstances, impose too much binding on the relative inclinations between the distinctive options.

This point is particularly vital to consider if the investigation expects to foresee how decisions would change in the event that one option was to vanish (for example on the off-chance that one political hopeful pulls back from a three applicant race). Different variations of the logistic model might be utilized as a part of such cases as they take into account the infringement of the IIA.

Probit regression

The Probit method is similar to the logistic method with the difference that these methods use different distributions. The logistic method uses logistic distribution, while the probit model uses z-distribution.

The z-distribution (used for probit) is as follows:

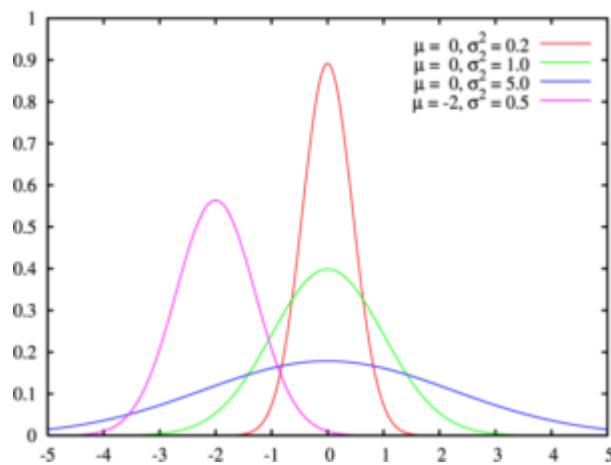


Image: "Normal distribution pdf" by D.328. License: [CC BY-SA 3.0](https://creativecommons.org/licenses/by-sa/3.0/)

Tools and Techniques Used When Only the Independent Variable Is Categorical

Simple regression

Suppose the independent variable of a particular study is gender (i.e., a categorical variable) and the dependent variable is earnings (i.e., a continuous variable), we could

run a simple regression to analyze the impact of gender (i.e., a categorical variable) on earnings (i.e., a continuous variable).

Suppose that the categorical variable 'gender' is coded in a manner that, if gender is 'male,' then it is coded as 1 and, if the gender is 'female,' then it is coded as 0. If the coefficient turns out to be 63.48, then it would imply that, if the person is a male, then the earnings would be 63.48 \$ higher than if the person was a female.

References

Chapnerkar, V. D. (n.d.). Live Stress-Free with Statistics and Numbers.

Heiman. (2013). Basic Statistics for the Behavioral Sciences. Cengage Learning.

Jayalakshmi, D. P., Pranitha, K., & Jyoti, Y. S. (2016). [Application of Mathematics in Economic Analysis](http://www.bomsr.com/4.S1.16/F13%20-86-92.pdf). Bulletin of Mathematics and Statistics Research, 4. Available at: <http://www.bomsr.com/4.S1.16/F13%20-86-92.pdf>

Legal Note: Unless otherwise stated, all rights reserved by Lecturio GmbH. For further legal regulations see our [legal information page](#).

Notes