

Welcome back. You are signed into
your member account
prashant9501@gmail.com. Not you?

✦ Member-only story

Chi-Square Test, with Python

The Complete Beginner's Guide to perform Chi-Square Test (with code!)



Chao De-Yu · [Follow](#)

Published in Towards Data Science

4 min read · May 22, 2021

Listen

Share

More



Photo by [Kalen Emsley](#) on [Unsplash](#)

In this article, I will introduce the fundamental of the chi-square test (χ^2), a statistical method to make the inference about the distribution of a variable or to decide whether there is a relationship exists between two variables of a population.

The inference relies on the χ^2 distribution curve dependent upon the number of degrees of freedom d.f.

Welcome back. You are signed into
your member account
prashant9501@gmail.com.

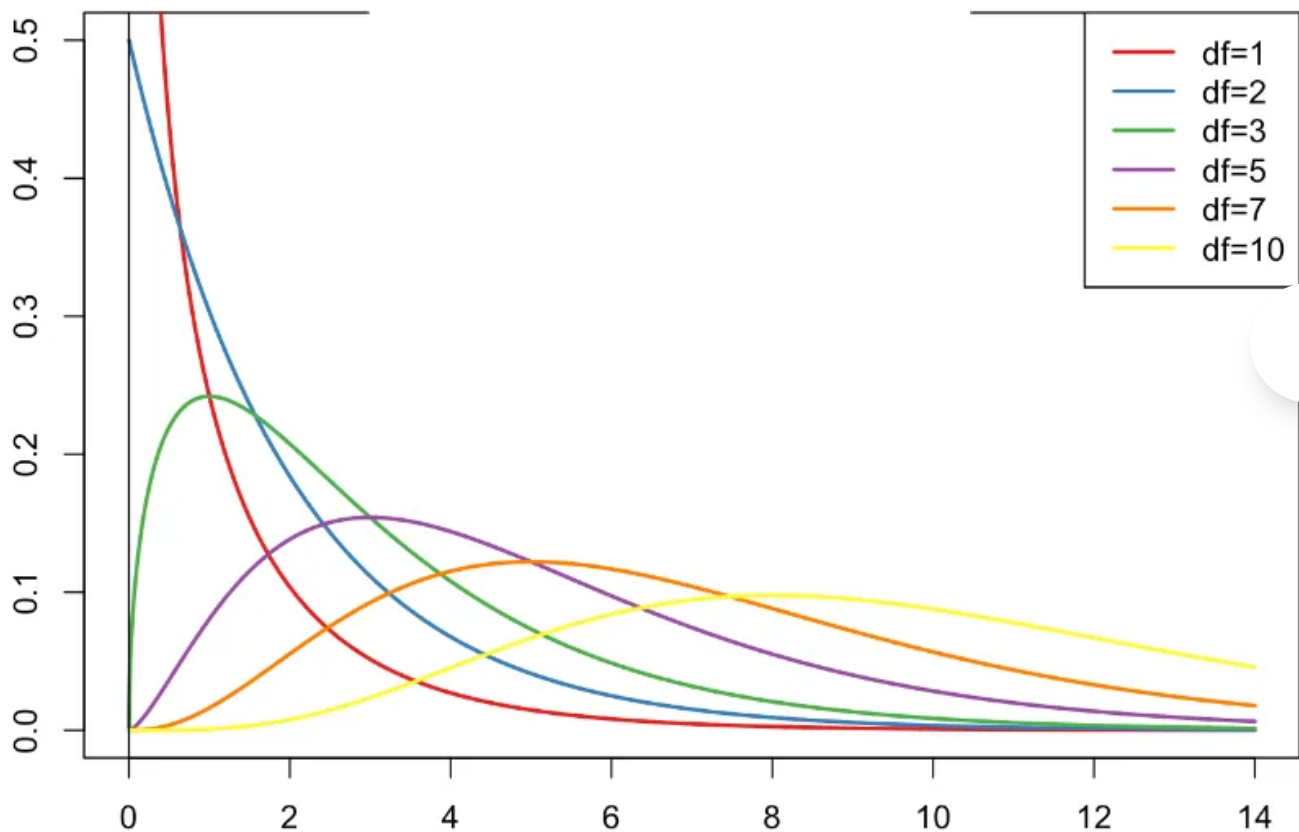


Figure 1: Chi-square distribution with different degree of freedom [1]

The χ^2 distribution curve is right-skewed and as the number of degrees of freedom becomes larger, the χ^2 curve will more similar to the normal distribution.

A: χ^2 test of Independence

It is used to decide whether there is a relationship exists between two variables of a population. Useful when analyzing survey results of 2 categorical variables.

- H_0 : The two categorical variables have **no relationship**
 H_1 : There is a **relationship** between two categorical variables
- The number of degrees of freedom of the χ^2 independence test statistics:
 $d.f. = (\# \text{ rows} - 1) * (\# \text{ columns} - 1)$

X_1 = row categorical variable with r levels X_2 = column categorical variable with c levels

		Welcome back. You are signed into your member account prashant9501@gmail.com.			
Row Variable X_1	C				Row Totals
R_1	C_1	O_{11}	O_{12}	O_{1c}	R_1 Total
R_2	C_1	O_{21}	O_{22}	O_{2c}	R_2 Total
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
R_r	C_1	O_{r1}	O_{r2}	O_{rc}	R_r Total
Column Totals	C_1 Total	C_2 Total	...	C_c Total	Grand Total

Table 1: rxc Contingency Table for 2 Categorical Variable

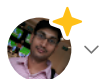
- If H_0 is true, each cell of the value in the contingency table above will contain a theoretical or expected frequency E_{ij} , as opposed to the observed value O_{ij} for each cell.

Assumed independent:

$$\begin{aligned}
 E_{ij} &= (R_i \cap C_j) = P(R_i) \times P(C_j) \times \text{Grand Total} \\
 &= \frac{R_i \text{ Total}}{\text{Grand Total}} \times \frac{C_j \text{ Total}}{\text{Grand Total}} \times \text{Grand Total} \\
 &= \frac{R_i \text{ Total} \times C_j \text{ Total}}{\text{Grand Total}}
 \end{aligned}$$

Figure 2: Derivation of the expected frequency

Open in app ↗



$$\chi^2 = \sum \frac{(O - E)^2}{E} \sim \chi^2_{(r-1)(c-1)}$$

over ALL the cells in the rxc contingency table

The table below is an exit poll survey of people in categories from their political affiliation, which is "Conservative, Socialist and Other". Is there any evidence of a relationship between the age group and their political affiliation, at 5% significant level?

	Conservative	Socialist	Other	Total
18-29	141	68	4	213
30-44	179	159	7	345
45-64	220	216	4	440
65 & older	86	101	4	191
Total	626	544	19	1189

Table 2: Exit poll survey [2]

According to five steps process of hypothesis testing:

H_0 : whether age group and their political affiliation are independent, i.e. no relationship

H_1 : whether age group and their political affiliation are dependent, i.e. \exists a relationship

$\alpha = 0.05$

Following χ^2 independence test statistics:

```

1  import pandas as pd
2  import scipy.stats as  Welcome back. You are signed into
3                               your member account
4  # create sample data : prashant9501@gmail.com.
5  data = [['18-29', 'Conservative'] for i in range(141)] + \
6          [['18-29', 'Socialist'] for i in range(68)] + \
7          [['18-29', 'Other'] for i in range(4)] + \
8          [['30-44', 'Conservative'] for i in range(179)] + \
9          [['30-44', 'Socialist'] for i in range(159)] + \
10         [['30-44', 'Other'] for i in range(7)] + \
11         [['45-65', 'Conservative'] for i in range(220)] + \
12         [['45-65', 'Socialist'] for i in range(216)] + \
13         [['45-65', 'Other'] for i in range(4)] + \
14         [['65 & older', 'Conservative'] for i in range(86)] + \
15         [['65 & older', 'Socialist'] for i in range(101)] + \
16         [['65 & older', 'Other'] for i in range(4)]
17  df = pd.DataFrame(data, columns = ['Age Group', 'Political Affiliation'])
18
19  # create contingency table
20  data_crosstab = pd.crosstab(df['Age Group'],
21                              df['Political Affiliation'],
22                              margins=True, margins_name="Total")
23
24  # significance level
25  alpha = 0.05
26
27  # Calculation of Chisquare
28  chi_square = 0
29  rows = df['Age Group'].unique()
30  columns = df['Political Affiliation'].unique()
31  for i in columns:
32      for j in rows:
33          O = data_crosstab[i][j]
34          E = data_crosstab[i]['Total'] * data_crosstab['Total'][j] / data_crosstab['Total']['Total']
35          chi_square += (O-E)**2/E
36
37  # The p-value approach
38  print("Approach 1: The p-value approach to hypothesis testing in the decision rule")
39  p_value = 1 - stats.chi2.cdf(chi_square, (len(rows)-1)*(len(columns)-1))
40  conclusion = "Failed to reject the null hypothesis."
41  if p_value <= alpha:
42      conclusion = "Null Hypothesis is rejected."
43
44  print("chisquare-score is:", chi_square, " and p value is:", p_value)
45  print(conclusion)
46
47  # The critical value approach
48  print("\n")

```



```

48 print("-----")
49 print("Approach 2: The critical value approach to hypothesis testing in the decision rule")
50 critical_value = stats.chi2.ppf(1 - alpha, df)
51 conclusion = "Failed to reject the null hypothesis. Welcome back. You are signed into your member account."
52 if chi_square > critical_value:
53     conclusion = "Null Hypothesis is rejected."
54
55 print("chisquare-score is:", chi_square, " and critical value is:", critical_value)
56 print(conclusion)

```

Approach 1: The p-value approach to hypothesis testing in the decision rule
chisquare-score is: 24.367421717305202 and p value is: 0.0004469083391495099
Null Hypothesis is rejected.

Approach 2: The critical value approach to hypothesis testing in the decision rule
chisquare-score is: 24.367421717305202 and critical value is: 12.591587243743977
Null Hypothesis is rejected.

Conclusion: We have enough evidence that there is an association between age group and their political affiliation, at 5% significance level.

B: χ^2 Goodness-Of-Fit Test

It is used to make the inference about the distribution of a variable.

- H_0 : The variable has the specified distribution, normal
 H_1 : The variable does not have the specified distribution, not normal
- The number of degrees of freedom of the χ^2 Goodness-Of-Fit test statistics:
d.f. = (# categories - 1)
- It compares the **observed frequencies O** of a sample with the **expected frequencies E**.
E = probability of the event * total sample size

The table below displays the more than 44 million people voting result for 2013 German Federal Election. 41.5% of German vote for the Christian Democratic Union (CDU), 25.7% for the Social Democratic Party (SPD) and the remaining 32.8% as Others.

Assume the researcher take a random sample and pick 123 students of FU Berlin about their party affiliation. These number correspond to the following table:

Welcome back. You are signed into your member account
prashant9501@gmail.com.

PD and 40 for Others.

Party	Percentage	Relative frequency
CDU	41.5	0.415
SPD	25.7	0.257
Others	32.8	0.328
	100	1

Table 3: 2013 German Federal Election [3]

According to five steps process of hypothesis testing:

H_0 : The variable has the specified distribution, i.e. the observed and expected frequencies are roughly equal

H_1 : The variable does not have the specified distribution, not normal

$\alpha = 0.05$

Following χ^2 Goodness-Of-Fit test statistics:

```

1  # Creation of data
2  data = [['CDU', 0.415, 'Welcome back. You are signed into your member account', 0]]
3  df = pd.DataFrame(data, columns=['id', 'name', 'message', 'd_freq'])
4  df['expected_freq'] = prashant9501@gmail.com.
5
6  # significance level
7  alpha = 0.05
8
9  # Calculation of Chisquare
10 chi_square = 0
11 for i in range(len(df)):
12     O = df.loc[i, 'observed_freq']
13     E = df.loc[i, 'expected_freq']
14     chi_square += (O-E)**2/E
15
16 # The p-value approach
17 print("Approach 1: The p-value approach to hypothesis testing in the decision rule")
18 p_value = 1 - stats.chi2.cdf(chi_square, df['Varname'].nunique() - 1)
19 conclusion = "Failed to reject the null hypothesis."
20 if p_value <= alpha:
21     conclusion = "Null Hypothesis is rejected."
22
23 print("chisquare-score is:", chi_square, " and p value is:", p_value)
24 print(conclusion)
25
26 # The critical value approach
27 print("\n-----")
28 print("Approach 2: The critical value approach to hypothesis testing in the decision rule")
29 critical_value = stats.chi2.ppf(1-alpha, df['Varname'].nunique() - 1)
30 conclusion = "Failed to reject the null hypothesis."
31 if chi_square > critical_value:
32     conclusion = "Null Hypothesis is rejected."
33
34 print("chisquare-score is:", chi_square, " and critical value is:", critical_value)
35 print(conclusion)

```

chisquare goodness of fit test by hosted with ❤️ by GitHub

view raw


Approach 1: The p-value approach to hypothesis testing in the decision rule
chisquare-score is: 1.693614940576721 and p value is: 0.42878164729702506
Failed to reject the null hypothesis.

Approach 2: The critical value approach to hypothesis testing in the decision rule
chisquare-score is: 1.693614940576721 and critical value is: 5.991464547107979
Failed to reject the null hypothesis.

Conclusion: We do not have enough evidence that the observed and expected frequencies are not equal

Welcome back. You are signed into your member account
prashant9501@gmail.com.

Recommended Reading

<div>ANOVA Test, with Python</div> <div>The Complete Beginner's Guide to perform ANOVA Test (with code!)</div> <div>towardsdatascience.com</div>	
<div>Two-Way ANOVA Test, with Python</div> <div>The Complete Beginner's Guide to perform Two-Way ANOVA Test (with code!)</div> <div>towardsdatascience.com</div>	
<div>McNemar's Test, with Python</div> <div>The Complete Beginner's Guide to perform McNemar's Test (with code!)</div> <div>towardsdatascience.com</div>	
<div>One-Sample Hypothesis Tests, with Python</div> <div>The Complete Beginner Guide to perform One-Sample Hypothesis Tests (with code!)</div> <div>levelup.gitconnected.com</div>	
<div>Two-Sample Hypothesis Tests, with Python</div>	

The Complete Beginner Guide to perform Two-Sample Hypothesis Tests (with code!)

levelup.gitconnected.com

Welcome back. You are signed into
your member account
prashant9501@gmail.com.

References

[1] “Chi-Square Tests • SOGA • Department of Earth Sciences.” [Online]. Available: <https://www.geo.fu-berlin.de/en/v/soga/Basics-of-statistics/Hypothesis-Tests/Chi-Square-Tests/index.html>

[2] “The Chi-Square Independence Test • SOGA • Department of Earth Sciences.” [Online]. Available: <https://www.geo.fu-berlin.de/en/v/soga/Basics-of-statistics/Hypothesis-Tests/Chi-Square-Tests/Chi-Square-Independence-Test/index.html>

[3] “Chi-Square Goodness-of-Fit Test • SOGA • Department of Earth Sciences.” [Online]. Available: <https://www.geo.fu-berlin.de/en/v/soga/Basics-of-statistics/Hypothesis-Tests/Chi-Square-Tests/Chi-Square-Goodness-of-Fit-Test/index.html>

Chi Square Test

Statistics

Data Science

Data Analysis



Follow

Written by Chao De-Yu

318 Followers · Writer for Towards Data Science

Data Analyst | MSc. Artificial Intelligence | LinkedIn — <https://www.linkedin.com/in/thet-thet-yee-deyu/>