

Chi Square test and contingency table

χ^2 Test

Another test that can be used for collocation detection.

5.3.3 Pearson's chi-square test

Use of the t test has been criticized because it assumes that probabilities are approximately normally distributed, which is not true in general (Church and Mercer 1993: 20). An alternative test for dependence which does not assume normally distributed probabilities is the χ^2 test (pronounced 'chi-square test'). In the simplest case, the χ^2 test is applied to 2-by-2 tables like table 5.8. The essence of the test is to compare the observed frequencies in the table with the frequencies expected for independence. If the difference between observed and expected frequencies is large, then we can reject the null hypothesis of independence.

Table 5.8 shows the distribution of *new* and *companies* in the reference corpus that we introduced earlier. Recall that $C(\text{new}) = 15,828$, $C(\text{companies}) = 4,675$, $C(\text{new companies}) = 8$, and that there are 14,307,668 tokens in the corpus. That means that the number of bi-grams $w_i w_{i+1}$ with the first token not being *new* and the second token being *companies* is $4667 = 4675 - 8$. The two cells in the bottom row are computed in a similar way.

The χ^2 statistic sums the differences between observed and expected values in all squares of the table, scaled by the magnitude of the expected values, as follows:

$$(5.6) \quad \chi^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

χ^2 Test

Für den χ^2 -Test brauchen wir die Kontingenztafel:

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} = 15820$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} = 4667$	$O_{22} = 14287173$	$O_{2-} = 14291840$
	$O_{-1} = 4675$	$O_{-2} = 14302993$	$O_{--} = 14307668$

Die χ^2 -Teststatistik wird folgendermaßen berechnet:

$$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad \text{wobei } E_{ij} = p_{i-} p_{-j} O_{--} = \frac{O_{i-} O_{-j}}{O_{--}}$$

Der χ^2 -Wert misst die Abweichung von den erwarteten Werten in der Kontingenz-Tabelle. Je größer er ist, desto kleiner ist der p-Wert.

Im Beispiel erhalten wir einen χ^2 -Wert von 1.55, der einem p-Wert von 0.21 entspricht. Das Ergebnis ist also auch bei diesem Test **nicht signifikant**.

Wie die Ergebnisse zeigen, können sich die p-Werte verschiedener statistischer Tests deutlich unterscheiden.

Method:

- fill the Contingency table
- compute X^2
- compute p-value of X^2
- see if the p-value is significant (it is significant if it is lower than the significant level that we picked)
- decide if we will reject the null hypothesis or not

<https://www.socscistatistics.com/pvalues/chidistribution.aspx>

Report a Chi-Square Result (APA)

Chi-square score:

1.55

DF:

1

Significance Level:

☐ 0.01

☒ 0.05

☐ 0.10

The P-Value is .213135. The result is *not* significant at $p < .05$.

Calculate

χ^2 Test

Für den χ^2 -Test brauchen wir die Kontingenztafel:

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} = 15820$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} = 4667$	$O_{22} = 14287173$	$O_{2-} = 14291840$
	$O_{-1} = 4675$	$O_{-2} = 14302993$	$O_{--} = 14307668$

Die χ^2 -Teststatistik wird folgendermaßen berechnet:

$$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \quad \text{wobei } E_{ij} = p_{i-} p_{-j} O_{--} = \frac{O_{i-} O_{-j}}{O_{--}}$$

O_{ij} ist der beobachtete Wert aus der Kontingenztafel.

E_{ij} sind die erwarteten Werte unter Annahme der Nullhypothese.

$p_{1-} = \frac{O_{1-}}{O_{--}}$ ist die Wahrscheinlichkeit, dass *new* das 1. Wort ist.

$p_{2-} = \frac{O_{2-}}{O_{--}}$ ist die Wahrscheinlichkeit, dass *new* **nicht** das 1. Wort ist.

$p_{-1} = \frac{O_{-1}}{O_{--}}$ ist die Wahrscheinlichkeit, dass *companies* das 2. Wort ist.

$p_{-2} = \frac{O_{-2}}{O_{--}}$ ist die Wahrscheinlichkeit, dass *companies* **nicht** das 2. Wort ist.

How to fill the table

set of words that are not „companies“

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} =$	$O_{12} =$	$O_{1-} =$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$	$O_{2-} =$
	$O_{-1} =$	$O_{-2} =$	$O_{--} =$

the number of (W1=new, W2= companies)

the number of
(W1=any word that is **not** „new“, W2= companies)
e.g. (W1=the, W2= companies)

the number of
(W1= any word, W2= companies)
e.g. (W1=new, W2= companies)
(W2=the, W2= companies)

- Fill what we know
 - $f(\text{new}) = 15828$
 - $f(\text{companies}) = 4675$
 - $f(\text{new}, \text{company}) = 8$
 - $n = 14307668$

new: 15828
 n=14307668
 company: 4675

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} =$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$	$O_{2-} =$
	$O_{-1} = 4675$	$O_{-2} =$	$O_{--} = 14307668$

Fill the rest (in the correct order)

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} =$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$	$O_{2-} =$
	$O_{-1} = 4675$	$O_{-2} =$	$O_{--} = 14307668$

$$\begin{aligned} O_{-2} &= O_{--} - O_{-1} \\ &= 14307668 - 4675 \\ &= 14,302,993 \end{aligned}$$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} =$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$	$O_{2-} =$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$	$O_{--} = 14307668$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} =$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$

$O_{1-} = 15828$
$O_{2-} =$
$O_{--} = 14307668$

$$\begin{aligned}
 O_{2-} &= O_{--} - O_{-1} \\
 &= 14307668 - 15828 \\
 &= 14,291,840
 \end{aligned}$$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} =$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$
		$O_{--} = 14307668$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} =$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$	$O_{2-} = 14,291,840$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$	$O_{--} = 14307668$

$$\begin{aligned}
 O_{12} &= O_{1-} - O_{11} \\
 &= 15828 - 8 \\
 &= 15,820
 \end{aligned}$$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} = 15,820$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$	$O_{2-} = 14,291,840$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$	$O_{--} = 14307668$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} = 15,820$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} =$	$O_{22} =$	$O_{2-} = 14,291,840$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$	$O_{--} = 14307668$

$$4675 - 8 = 4667$$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} = 15,820$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} = 4667$	$O_{22} =$	$O_{2-} = 14,291,840$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$	$O_{--} = 14307668$

	$w_2 = \text{companies}$	$w_2 \neq \text{companies}$	
$w_1 = \text{new}$	$O_{11} = 8$	$O_{12} = 15,820$	$O_{1-} = 15828$
$w_1 \neq \text{new}$	$O_{21} = 4667$	$O_{22} = 14287173$	$O_{2-} = 14,291,840$
	$O_{-1} = 4675$	$O_{-2} = 14,302,993$	$O_{--} = 14307668$

$$14,302,993 - 15820 = 14287173$$

finished