

χ^2 Test

Für den χ^2 -Test verwenden wir die folgende 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:

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

O_{ij} ist der Wert aus der Kontingenztafel

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

Das Signifikanzniveau, das der Statistik entspricht, schlägt man in einer χ^2 -Statistiktafel nach. (Beispiel: $\chi^2 \approx 1.55 < 3.84 \Rightarrow$ nicht signif.)

A **chi-squared test**, also written as **χ^2 test**, is a **statistical hypothesis test** that is **valid** to perform when the test statistic is **chi-squared distributed** under the **null hypothesis**, specifically **Pearson's chi-squared test** and variants thereof. Pearson's chi-squared test is used to determine whether there is a **statistically significant** difference between the expected **frequencies** and the observed frequencies in one or more categories of a **contingency table**.

Main Idea: hypothesis testing

We want to test if "new companies" is a collocation or not.

Given observed information, the hypothesis testing result tells us whether "new companies" happens more than expected or not.

Our assumption is that, if the word pair happens more than expected then it could be because it is a collocation

-but we can not be so sure that the information we have (what we observe) is enough to make this claim.

-**hypothesis testing** can tell us if there is enough evidence (significant deviation from the expected value) that our hypothesis is correct or not

Pearson's chi-squared test is used to determine whether there is a statistically significant difference between the expected frequencies and the observed frequencies in one or more categories of a contingency table.

Beispiel: Ist das Wortpaar **new companies** signifikant häufiger als erwartet?

Daten: In einem Korpus mit $n=14,307,668$ Wörtern, taucht **new** $f_{new}=15,828$ Mal auf, **companies** $f_{companies}=4,675$ Mal und **new companies** 8 Mal.

χ^2 Test

Für den χ^2 -Test verwenden wir die folgende 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:

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

O_{ij} ist der Wert aus der Kontingenztafel

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

Das Signifikanzniveau, das der Statistik entspricht, schlägt man in einer χ^2 -Statistiktafel nach. (Beispiel: $\chi^2 \approx 1.55 < 3.84 \Rightarrow$ nicht signif.)

Explain the notations

O_position1, position2

* pos1 represents the existence w1. It can take 2 values. Value 1 means the word exists, 2 means does not exist.

* pos2 represents w2.

11: $f(w_1 = \text{new}, w_2 = \text{company})$

22: $f(w_1 \neq \text{new}, w_2 \neq \text{company})$

12: $f(w_1 = \text{new}, w_2 \neq \text{company})$

21: $f(w_1 \neq \text{new}, w_2 = \text{company})$

1_ : $f(w_1 = \text{new}, w_2 = \text{any})$

_1: $f(w_1 = \text{any}, w_2 = \text{company})$

_2: $f(w_1 = \text{any}, w_2 \neq \text{company})$

2_ : $f(w_1 \neq \text{new}, w_2 = \text{any})$

_ _ = : $f(w_1 = \text{any}, w_2 = \text{any})$

Example

-create a contingency table for testing “new companies”

-In a corpus of 100 words, freq(new) = 20, freq(companies) = 40, freq(new companies) = 10

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

Put the values from the table in the formula

Example

-create a contingency table for testing “polar bear”

-In a corpus of 100 words, $\text{freq}(\text{polar}) = 20$, $\text{freq}(\text{bear}) = 40$, $\text{freq}(\text{polar bear}) = 10$, fill the table

	w2	not w2	
w1	$O_{11} = 10 \text{ f}(\text{polar}, \text{bear}) \text{ (1)}$	$O_{12} = \text{f}(\text{polar}, \text{not bear}) \text{ (8)}$ $= \text{f}(\text{polar}, \text{any}) - \text{f}(\text{polar}, \text{bear})$ $O_{12} = O_{1_} - O_{11}$	$O_{1_} = 20 \text{ f}(\text{polar}, \text{any}) \text{ (3)}$
not w1	$O_{21} = \text{f}(\text{not polar}, \text{bear}) \text{ (7)}$ $= \text{f}(\text{any}, \text{bear}) - \text{f}(\text{polar}, \text{bear})$ $O_{21} = O_{_1} - O_{11}$	$O_{22} = \text{f}(\text{not polar}, \text{not bear})$ $= \text{f}(\text{any}, \text{not bear}) - \text{f}(\text{polar}, \text{not bear}) \text{ (9)}$ $O_{22} = O_{_2} - O_{12}$	$O_{2_} = \text{f}(\text{not polar}, \text{any}) \text{ (6)}$ $= \text{f}(\text{any}, \text{any}) - \text{f}(\text{polar}, \text{any})$ $O_{2_} = N - O_{1_}$
	$O_{_1} = 40 \text{ f}(\text{any}, \text{bear}) \text{ (2)}$	$O_{_2} = \text{f}(\text{any}, \text{not bear}) \text{ (5)}$ $= \text{f}(\text{any}, \text{any}) - \text{f}(\text{any}, \text{bear})$ $O_{_2} = N - O_{_1}$	$O_{_} = 100 \text{ (4) f}(\text{any}, \text{any})$

Trick on how to calculate the score

- You will be given 4 values
- Start by writing the values in the correct cells (cells with number 1, 2, 3, 4)
- see the table as two parts: inner part (blue cells) and out part (light green cells)
- next, you will fill cell 5, 6, 7, 8, 9
- to compute cell 5, you calculate cell 4 - cell 2, $[4] [5] [2]$
- then to compute cell 6, you calculate cell 4 - cell 3, $[4] [6] [3]$
- now the outer cells are computed, next you want to compute the inner cells
- you start by calculating cell 7 by taking cell 2 - cell 1
- then to compute cell 8, you take cell 3 - cell 1, $[1] [8] [3]$
- for the last cell 9, you take cell 5 - cell 8, $[8] [9] [5]$
- You will see a pattern here (if you have 3 cell $[a] [] [c]$, you compute the middle cell by taking $c - a$)
- $[a]$
- $[]$
- $[c]$
- for this you take $c - a$ as well (the outer cell c has higher value, so we subtract a from c)
- start from filling most outer cells, then going to the inner cells

Übung:

-create a contingency table for testing “blue eyes”

-In a corpus of 1000 words, $\text{freq}(\text{blue}) = 25$, $\text{freq}(\text{eyes}) = 30$, $\text{freq}(\text{blue eyes}) = 15$