

Signifikanzmaße bei der Beurteilung von Kookkurrenzen

In der Statistik wird unter Signifikanz eine Kennzahl verstanden, welche die Wahrscheinlichkeit eines systematischen Zusammenhangs zwischen Variablen, im Falle von Textanalysen also zwischen Teiltextrn (z.B. Wörtern), bezeichnet. Die Signifikanz drückt aus, ob ein scheinbarer Zusammenhang rein zufälliger Natur sein könnte oder mit hoher Wahrscheinlichkeit tatsächlich vorliegt.

Zur Berechnung werden abhängig vom Untersuchungsgegenstand unterschiedliche Formeln herangezogen, welche in erster Linie aus der Computerlinguistik stammen. Die Signifikanzmaße sollen dabei helfen, wichtige von unwichtigen Kookkurrenzen zu trennen. Dabei werden statistische Kenngrößen, wie Korpusgröße, Häufigkeit der einzelnen Wörter oder Frequenz des gemeinsamen Auftretens, ins Verhältnis gesetzt.

Eines der einfachsten Signifikanzmaße ist eine frequenzsortierte Kookkurrenzliste, also die Häufigkeit des gemeinsamen Auftretens zweier Worte im Gesamtkorpus. Ein Nachteil frequenzsortierter Listen ist, dass nach dem Zipf'schen Gesetz, dem Beginn der quantitativen Linguistik, sehr viele Wörter sehr selten auftreten. Demzufolge lassen sich mit einem Schwellenwert größer 1, also dem mehrmaligen gemeinsamen Auftreten eines Wortpaares, etwa zwei Drittel der Kookkurrenzen herausfiltern. Berechnet von den eAQUA-Tools sieht dies für ausgewählte Korpora wie folgt aus:

Korpus	Anzahl Kookkurrenzen	Kookkurrenzen freq = 1	in Prozent
BTL ¹⁾	137.486.214	110.876.836	80,65
MPL ²⁾	580.247.568	398.935.822	68,75
Perseus Shakespeare ³⁾	6.746.602	5.027.170	74,51
TLG ⁴⁾	355.021.014	258.961.566	72,94

Wie aus der kleinen Übersicht zu erkennen ist, sind ein Großteil der gefundenen Kookkurrenzen eher als niedrigfrequent zu bezeichnen. Um daraus die wichtigen zu filtern, sind Berechnungsmethoden erforderlich, von denen hier einige vorgestellt werden.

Dice

Beim Dice-Koeffizienten (auch Sørensen-Dice-Koeffizient, benannt nach den Botanikern Thorvald Sørensen und Lee Raymond Dice) wird die Ähnlichkeit zweier Terme mittels einer Zahl zwischen 0 und 1 angegeben. Berechnungsgrundlage sind sogenannte N-Gramme. Bei N-Grammen wird ein Term bzw. ein Text in gleich große Teile zerlegt. Diese Fragmente können Buchstaben, Phoneme, ganze Wörter oder ähnliches sein.

Ermittelt wird die Anzahl der N-Gramme, die in beiden Termen vorhanden sind, um diese ins Verhältnis zur

Gesamtzahl der N-Gramme zu setzen. Berechnet wird nach der Formel
$$dice_{ab} = \frac{2 \times n_{ab}}{n_a + n_b}$$
 wobei n_{ab} die Schnittmenge beider Terme und n_a bzw. n_b die Anzahl der gebildeten N-Gramme pro Term angibt.

Beispiel 1: Ausdruck a = Tür Ausdruck b = Tor	$dice_{ab} = \frac{2 \times n_{ab}}{n_a + n_b}$
Bigramm	Trigramm
$a = \{ \text{\$T, Tü, ür, r\$} \}$ $b = \{ \text{\$T, To, or, r\$} \}$ $d_{Tür, Tor} = \frac{2 \times 2}{4 + 4} = \frac{4}{8} = 0,5$	$a = \{ \text{\$\$T, \$Tü, TÜR, ür$, r\$\$} \}$ $b = \{ \text{\$\$T, \$To, Tor, or$, r\$\$} \}$ $d_{Tür, Tor} = \frac{2 \times 2}{5 + 5} = \frac{4}{10} = 0,4$

Beispiel 2 Ausdruck a = Spiegel Ausdruck b = Spargel	$dice_{ab} = \frac{2 \times n_{ab}}{n_a + n_b}$
Bigramm	Trigramm
a = { §S, Sp, pi, ie, eg, ge, el, l§ } b = { §S, Sp, pa, ar, rg, ge, el, l§ }	a = { §§S, §Sp, Spi, pie, ieg, ege, gel, el§, l§§ } b = { §§S, §Sp, Spa, par, arg, rge, gel, el§, l§§ }
$d_{Spiegel, Spargel} = \frac{2 \times 5}{8+8} = \frac{10}{16} = 0,625$	$d_{Spiegel, Spargel} = \frac{2 \times 5}{9+9} = \frac{10}{18} \approx 0,556$

Bei der Bewertung von Kookkurrenzen kann der Dice-Koeffizient genutzt werden, indem die Häufigkeiten (Frequenzen) der Wörter ins Verhältnis gesetzt werden. n_a und n_b sind dabei die Frequenzen der Terme, n_{ab} die Anzahl des gemeinsamen Auftretens.

Aus der oben angeführten Berechnung ergeben sich relativ einfache Bewertungsmaßstäbe. Je frequenter die beiden Begriffe gemeinsam benutzt werden, um so mehr nähert sich der Wert 1. Treten beide Begriffe nur gemeinsam auf, wird die höchste Signifikanz mit 1 erreicht. Wie oft diese Kookkurrenz im Korpus zu finden ist, spielt dabei keine Rolle. Daraus ergibt sich eine wichtige Eigenschaft des Dice-Koeffizienten: Kookkurrenzen, die selten zusammen auftreten, bei denen ein Wort hoch- und das andere niedrigfrequent sind, werden als unsignifikant bewertet.

Jaccard

Beim Jaccard-Koeffizienten (nach dem Botaniker Paul Jaccard) wird die Ähnlichkeit zweier Terme mittels einer Zahl zwischen 0 und 1 angegeben. Berechnungsgrundlage bei Textmining-Verfahren sind sogenannte N-Gramme. Bei N-Grammen wird ein Term bzw. ein Text in gleich große Teile zerlegt. Diese Fragmente können Buchstaben, Phoneme, ganze Wörter oder ähnliches sein.

Ermittelt wird die Anzahl der N-Gramme, die in beiden Termen vorhanden sind, um diese ins Verhältnis zur

Gesamtzahl der N-Gramme zu setzen. Berechnet wird nach der Formel $jaccard_{ab} = \frac{n_{ab}}{n_a + n_b - n_{ab}}$ wobei n_{ab} die Schnittmenge beider Terme und n_a bzw. n_b die Anzahl der gebildeten N-Gramme pro Term angibt.

Beispiel 1: Ausdruck a = Tür Ausdruck b = Tor	$jaccard_{ab} = \frac{n_{ab}}{n_a + n_b - n_{ab}}$
Bigramm	Trigramm
a = { §T, Tü, ür, r§ } b = { §T, To, or, r§ }	a = { §§T, §Tü, Tür, ür§, r§§ } b = { §§T, §To, Tor, or§, r§§ }
$d_{Tür, Tor} = \frac{2}{4+4-2} = \frac{2}{6} \approx 0,334$	$d_{Tür, Tor} = \frac{2}{5+5-2} = \frac{2}{8} = 0,25$
Beispiel 2 Ausdruck a = Spiegel Ausdruck b = Spargel	$jaccard_{ab} = \frac{n_{ab}}{n_a + n_b - n_{ab}}$
Bigramm	Trigramm
a = { §S, Sp, pi, ie, eg, ge, el, l§ } b = { §S, Sp, pa, ar, rg, ge, el, l§ }	a = { §§S, §Sp, Spi, pie, ieg, ege, gel, el§, l§§ } b = { §§S, §Sp, Spa, par, arg, rge, gel, el§, l§§ }
$d_{Spiegel, Spargel} = \frac{5}{8+8-5} = \frac{5}{11} \approx 0,455$	$d_{Spiegel, Spargel} = \frac{5}{9+9-5} = \frac{5}{13} \approx 0,385$

Für die Bewertung von Kookkurrenzen gilt beim Jaccard-Koeffizienten ähnliches wie beim Dice-Koeffizienten. Beide berechnen den Signifikanzwert ähnlich, die relative Ordnung der Kookkurrenzen bleibt gleich, nur der

absolute Signifikanzwert unterscheidet sich marginal. Eine Modell-Berechnung mit mittlerer Frequenz von 100 sieht wie folgt aus:

n_a	n_b	n_{ab}	Dice	Jaccard
100	100	1	0,01	0,005
100	100	10	0,1	0,05
100	100	50	0,5	0,33
100	100	90	0,9	0,82
100	100	100	1	1

Poisson-Maß

Ein Ansatz zur Berechnung von signifikanten Kookkurrenzen basiert auf der Poisson-Verteilung (benannt nach dem Mathematiker Siméon Denis Poisson), einer diskreten Wahrscheinlichkeitsverteilung $p(n, k) = \frac{1}{k!} \gamma^k e^{-\gamma}$

Auf der Basis der Poisson-Verteilung geben Quasthoff / Wolff⁵⁾ das Poisson-Maß mit der Formel

$$p(n_a, n_b, k, n) = \frac{k \times (\log k - \log \gamma - 1)}{\log n}$$

an, welche beispielsweise für die Berechnung von Korpora im Wortschatz-Portal genutzt wurde, und in der die zwei Faktoren n (Anzahl der Sätze im Korpus) und k (Häufigkeit des gemeinsamen Auftretens, auch n_{ab} bezeichnet) maßgeblich sind.

Nach einer Umstellung und der Grundannahme $\gamma = \frac{n_a \times n_b}{n}$ ergibt sich folgende Berechnung

$$p = \frac{n_{ab} \times \log \frac{n_a \times n_b}{n} - n_{ab}}{\log n}$$

Somit ließe sich das Poisson-Maß auf die Differenz zwischen Local Mutual Information und Frequenz reduzieren.

Log-Likelihood-Maß

Eine der populärsten Signifikanzmaße bei der Analyse großer Textcorpora ist nach Dunning⁶⁾ das Log-Likelihood-Maß, welches auf der Binomialverteilung, eine der wichtigsten diskreten Wahrscheinlichkeitsverteilungen, basiert.

$$p(K=k) = p^k (1-p)^{n-k} \binom{n}{k}$$

Dunning kommt schließlich bei der Berechnung von **log likelihood** zu der Formel:

$$-2 \log \lambda = 2 \left[\log L(p_1, k_1, n_1) + \log L(p_2, k_2, n_2) - \log L(p_1, k_1, n_1) - \log L(p_2, k_2, n_2) \right]$$

unter der Voraussetzung

$$\log L(p, n, k) = k \log p + (n - k) \log(1 - p)$$

Das Log-Likelihood-Maß kann demzufolge abgeleitet werden

$$|g| = 2 \left[n \log n - n_a \log n_a - n_b \log n_b + n_{ab} \log n_{ab} + \binom{n - n_a - n_b + n_{ab}}{n_a - n_b + n_{ab}} \log \binom{n - n_a - n_b + n_{ab}}{n_a - n_b + n_{ab}} + \binom{n - n_{ab}}{n_a - n_{ab}} \log \binom{n - n_{ab}}{n_a - n_{ab}} + \binom{n - n_{ab}}{n_b - n_{ab}} \log \binom{n - n_{ab}}{n_b - n_{ab}} - \binom{n - n_a}{n_b - n_a} \log \binom{n - n_a}{n_b - n_a} - \binom{n - n_b}{n_a - n_b} \log \binom{n - n_b}{n_a - n_b} \right]$$

Charakteristisch für das Log-Likelihood-Maß ist, im Gegensatz beispielsweise zum Poisson-Maß, die Gleichbehandlung von signifikant häufigen und signifikant seltenen Ereignissen. So finden sich in den Digitalisaten vom TLG in der Version TLG-E bei rund 73,8 Millionen Wörtern etwa 1,3 Millionen Kookkurrenzen, die nur einmal auftreten und trotzdem mit einem lgl-Wert von 30 und ein wenig mehr belegt sind. Einen ähnlich großen Wert von 34,553 haben zum Beispiel **καί** und **Τὸ**, die zusammen 14311 Mal gezählt wurden.

¹⁾

Bibliotheca Teubneriana Latina, Online-Version, Stand vom Februar 2014

²⁾

Patrologia Latina Database, CD-ROM Version, November 1995c

³⁾

William Shakespeare in Perseus Digital Library, Renaissance Materials, Stand vom Mai 2013

⁴⁾

TLG-E, CD-ROM Version aus dem Jahre 1999

⁵⁾

[Quasthoff 02]. Uwe QUASTHOFF, Christian WOLFF. The Poisson Collocation Measure and its Applications. In Second International Workshop on Computational Approaches to Collocations, 2002.

⁶⁾

[Dunning 93]. Dunning, T. „Accurate Methods for the Statistics of Surprise and Coincidence“. In: Computational Linguistics 19, 1 (1993), 61-74.

From:

<http://www.eaqua.net/doku/> - **Wissensdatenbank**

Permanent link:

<http://www.eaqua.net/doku/doku.php/signifikanz?rev=1400571148>

Last update: **2018/05/15 11:30**