

ASYMMETRIE IN MEHRWORTAUSDRÜCKEN: „krasser Unterschied“ vs „himmelweiter Unterschied“

Yana Strakatova (Tübingen)

Kookkurrenzen

1. Kollokat + Basis:

großer Unterschied, fürstliches Gehalt, blinder Passagier, schwarzes Schaf

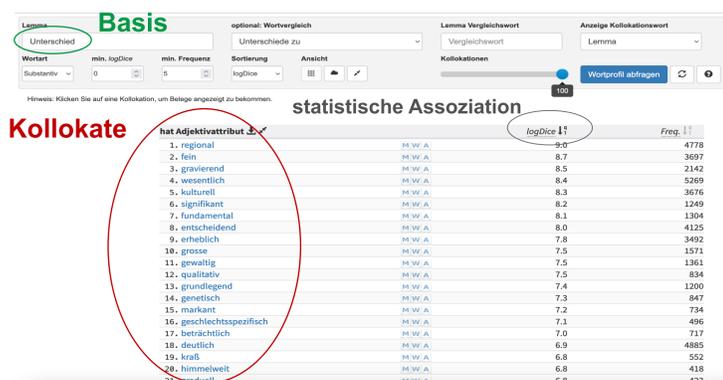
2. Statistische Assoziation:

Berechnet basierend auf Korpusfrequenzen („contingency table“ [1])

$$\text{Statistischer Assoziationswert } \logDice = 14 + \log_2 \frac{2a}{(a+c) + (a+b)} \quad [2]$$

- theoretische Obergrenze = 14, die meisten Werte < 10
- unabhängig von Korpusgröße
- negative Werte = keine statistische Signifikanz

DWDS-Wortprofil [3]



Kollokat (logDice)

- wesentlich (8.4)
- beträchtlich (7.0)
- deutlich (6.9)
- krass (6.8)
- himmelweit (6.8)

Contingency table

	Nomen	¬Nomen	Σ
Adjektiv	a	b	a+b
¬Adjektiv	c	d	c+d
Σ	a+c	b+d	a+b+c+d=n

a – Frequenz des Bigrams; a+c – Frequenz des Nomens;
a+b – Frequenz des Adjektives

Asymmetrie

krass: Unterschied, Gegensatz, Kontrast, Gegenteil, Ungleichheit, Diskrepanz, Abweichung, usw.
himmelweit: Unterschied

Asymmetrischer statistischer Assoziationswert: ΔP [4]
Werte zwischen -1 und 1

left-predictive: ΔP_1

$$\Delta P_1 = \Delta P(w_1|w_2) = \frac{a}{a+c} - \frac{b}{b+d}$$

← frühe Morgenstunde

right-predictive: ΔP_2

$$\Delta P_2 = \Delta P(w_2|w_1) = \frac{a}{a+b} - \frac{c}{c+d}$$

→ himmelweiter Unterschied

Wie häufig ist solche Asymmetrie? Welche Arten von Mehrwortausdrücken sind asymmetrisch? Kann man asymmetrische Werte in das DWDS-Wortprofil einbauen?

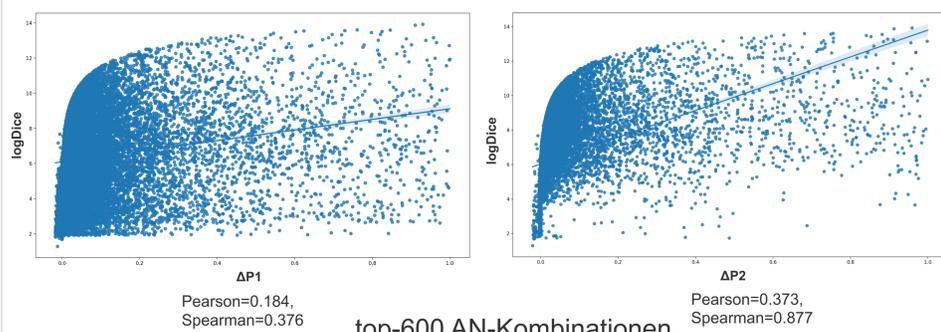
Korpusstudie

Daten

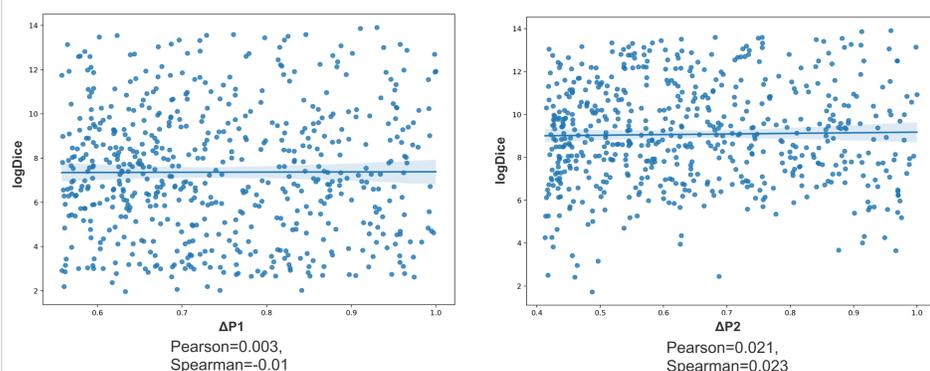
- Korpus (Wikipedia 2017, 2018 [5]; decow16ax [6])
- Frequenz > 400
- 36,945 AN-Kombinationen
- Statistische Werte: ΔP_1 , ΔP_2 , logDice
- Manuelle Annotation als „Kollokation“, „Begriff“, „Eigenname“, „Freie Phrase“, „Idiom“, „Klischee“

Korrelation mit logDice

alle AN-Kombinationen



top-600 AN-Kombinationen



600-best right-predictive (ΔP_2)

Label	Anzahl	Beispiele
Begriff	256	rheumatoide Arthritis (0.97), glykämischer Index (0.76), festverzinsliche Wertpapiere (0.64), molare Masse (0.54)
Kollokation	182	geraume Zeit (0.98), triftiger Grund (0.89), brenzlige Situation (0.84), faustdicke Überraschung (0.70), klirrende Kälte (0.57)
Eigenname	87	Timmendorfer Strand (0.93), (der) Gestiefelte Kater (0.90), Brühlsche Terrasse (0.73), (Das) Wohltemperierte Klavier (0.46)
Idiom	15	gordischer Knoten (0.89), zweischneidiges Schwert (0.61), (am) seidenen Faden (hängen) (0.41)

600-best left-predictive (ΔP_1)

Label	Anzahl	Beispiele
Eigenname	240	Olympische Sommerspiele (0.96), (der) Letzte Mohikaner (0.87), Hohe Tatra (0.62), Vereinigtes Königreich (0.56)
Begriff	208	vorläufige Vollstreckbarkeit (0.95), chronische Polyarthritis (0.83), absoluter Nullpunkt (0.68), Neue Sachlichkeit (0.65)
Freie Phrase	65	passender Staubsaugerbeutel (0.95), langer Blütenstiel (0.60), früher Sonntagmorgen (0.56)
Kollokation	40	letzte Ruhestätte (0.91), üble Nachrede (0.90), reine Formsache (0.89), gesunder Menschenverstand (0.87), stehende Ovation (0.79), schwere Kopfverletzung (0.61)
Idiom	28	innerer Schweinehund (0.85), eierlegende Wollmilchsau (0.85), ewige Jagdgründe (0.69), fliegende Untertasse (0.69)
Klischee	6	zu guter Letzt (0.91), gutes Gelingen (0.79), herzlichen Glückwunsch (0.78), gute Besserung (0.68)
Kompositum	1	rote Beete (0.63)