

N-gram Modelling

Katja Markert

Institut für Computerlinguistik
Uni Heidelberg
markert@cl.leeds.ac.uk

November 21, 2019

- Bis jetzt: Einzelne Wörter erkennen, vergleichen und zählen.
Wortverteilungen.
- Jetzt: Sprachmodelle: Verteilungen über Sätze, Texte,
Wortsequenzen
- n -gram Modelle: bestimmte Art von Sprachmodellen

- 1 N-gram Modelle: Grundlagen
- 2 Schätzung von n -gram Modellen und Smoothing
- 3 Praxisprobleme und Intrinsische Evaluation
- 4 Einfache Anwendungen: Confusion Sets und Sprachidentifikation

- 1 N-gram Modelle: Grundlagen
- 2 Schätzung von n -gram Modellen und Smoothing
- 3 Praxisprobleme und Intrinsische Evaluation
- 4 Einfache Anwendungen: Confusion Sets und Sprachidentifikation

Wir wollen herausfinden: Ist

$P(\text{He often wears tight trousers})$

wahrscheinlicher als?

$P(\text{He often wears tight socks})$

wahrscheinlicher als?

$P(\text{He often knits tight trousers})$

Wie wahrscheinlich ist es, dass *trousers* dem Anfang *He often wears tight* folgt?

- Spracherkennung
- Handschriftenerkennung
- Kontextsensitive Rechtschreibkorrektur
- Textvervollständigung: texting
- Beurteilung von generierten Sprach- oder Textsequenzen in Summarization, Machine Translation . . .

Viele Korrekturprogramme nutzen Wahrscheinlichkeiten von Zeichenketten, um ungewöhnliche Sequenzen zu erkennen und wahrscheinlichere Sequenzen vorherzusagen

Wahrscheinlichkeit $P(t)$ für Wortsequenz $t = w_1 w_2 \dots w_n$ nach Kettenregel:

$$\begin{aligned} P(t) = P(w_1 w_2 \dots w_n) &= P(w_1)P(w_2|w_1) \dots P(w_n|w_1, \dots, w_{n-1}) \\ &= \prod_{i=1}^n P(w_i|w_1 \dots w_{i-1}) \end{aligned}$$

$$P(w_i|w_1 \dots w_{i-1}) = \frac{P(w_1 \dots w_i)}{P(w_1 \dots w_{i-1})}$$

Gleichwertigkeit von Sequenzmodellierung und Vorhersage

Sehe $\langle s \rangle$ als spezielles Zeichen für Satzanfang und $\langle /s \rangle$ als spezielles Zeichen für Satzende

Damit

$$\begin{aligned}P(s) &= P(\langle s \rangle w_1 w_2 \dots w_n \langle /s \rangle) \\ &= P(\langle s \rangle)P(w_1 | \langle s \rangle)P(w_2 | \langle s \rangle w_1) \dots \\ &\quad P(w_n | \langle s \rangle w_1, \dots, w_{n-1}) \cdot P(\langle /s \rangle | \langle s \rangle w_1 \dots w_n)\end{aligned}$$

Markovannahme

w_i hängt nur von den Tokens unmittelbar davor ab.

- 0te Ordnung = Unigramme

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i)$$

- Erste Ordnung = Bigramme

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | w_{i-1}) = \frac{P(w_{i-1} w_i)}{P(w_{i-1})}$$

- n -te Ordnung parallel
- Meist benutzt (Zweite Ordnung, Trigramme)

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | w_{i-2} w_{i-1}) = \frac{P(w_{i-2} w_{i-1} w_i)}{P(w_{i-2} w_{i-1})}$$

Bigram Model

$$P(\langle s \rangle \text{ He often wears tight trousers } \langle /s \rangle) \approx \\ P(\langle s \rangle)P(\text{He} | \langle s \rangle)P(\text{often} | \text{He})P(\text{wears} | \text{often})P(\text{tight} | \text{wears}) \\ P(\text{trousers} | \text{tight}) \cdot P(\langle /s \rangle | \text{trousers})$$

► unigram

► trigram

Wenn Sie sich folgendes Beispiel anschauen, welche Vor- und Nachteile sehen Sie dann für n-gram Modelle?

He doesn't wear loose trousers, he prefers to

When I

He doesn't wear loose socks he prefers to

wear tight x

The **man** that I was watching without pausing to look at what was happening down the street, and quite oblivious to the situation that was about to befall him confidently **strode** into the center of the road.

- Das passiert nicht oft (im Englischen)
- Collins (1996): 74% der Abhängigkeiten in Penn Treebank gehören zu einem angrenzenden Wort

Collins, Michael John. *A new statistical parser based on bigram lexical dependencies*. Proceedings of the 34th annual meeting on Association for Computational Linguistics. Association for Computational Linguistics, 1996.

- Sprachmodelle sind Wahrscheinlichkeitsverteilungen über Sätze
- N-gram Modelle eine Art von Sprachmodellen
- basierend auf recht kurzer Historie
- In Anwendungen oft kombiniert mit anderen Restriktionen: wie häufigen Rechtschreibfehlern und Minimum Edit Distance
- * Jurafsky und Martin, Kapitel 3, 3rd/online edition
- Übungsblatt 5

- 1 N-gram Modelle: Grundlagen
- 2 Schätzung von n -gram Modellen und Smoothing**
- 3 Praxisprobleme und Intrinsische Evaluation
- 4 Einfache Anwendungen: Confusion Sets und Sprachidentifikation

Wie leite ich ein n -gram Modell in der Praxis her?

- Starte mit Vokabular V (word types):
 - alle types in einem Korpus;
 - oder die m häufigsten;
- Wähle Modelltyp:
 - Nullte Ordnung= unigram
 - Erste Ordnung = bigram
 - Zweite Ordnung = trigram
- Parameterschätzung e.g. $P(w_1, w_2)$

▶ Parameter Estimation

Relative Häufigkeit = **maximum likelihood Schätzung** auf Korpus der Grösse N .

$$P(w_1) = \frac{f(w_1)}{N}$$

$$P(w_2|w_1) = \frac{f(w_1, w_2)}{f(w_1)}$$

$$P(w_3|w_1, w_2) = \frac{f(w_1, w_2, w_3)}{f(w_1, w_2)}$$

Korpus:

< s > I am Sam < /s >

< s > Sam I am < /s >

< s > I do not like green eggs and ham < /s >

Teil eines Bigram-modells des Brown Korpus

Erstes Wort in Zeile.

	l	drank	white	wine	milk	rum	water	
l	0	4	0	0	0	0	0	3789
drank	0	0	0	0	0	0	0	19
white	0	0	0	5	1	0	0	360
wine	0	0	0	0	0	0	0	67
milk	0	0	0	0	0	0	0	46
rum	0	0	0	0	0	0	0	3
water	0	0	0	0	0	0	0	437

Brown corpus hat 1 m tokens.

Vergleiche auch mit Beispielen in Jurafsky und Martin, mit Satzanfängen

Teil eines Bigram-Modells auf dem BNC

Erstes Wort in Zeile

	l	drank	white	wine	milk	rum	water	
l	4492	93	0	0	0	0	6	868973
drank	1	0	1	6	3	4	16	1312
white	18	0	25	272	1	7	120	23445
wine	17	0	1	3	0	0	0	6061
milk	15	0	2	0	5	0	0	4737
rum	2	0	0	0	0	0	0	429
water	40	0	0	0	1	0	17	34349

BNC hat 100m Tokens.

Welche der folgenden Sequenzen sind wahrscheinlicher in einem Unigrammodell auf Brown/BNC? In einem Bigrammodell?

- 1 I drank white rum
- 2 I drank white water
- 3 I drank white wine
- 4 I drank white milk

Was passiert wohl in einem Trigram-Modell?

- Im Idealfall führt größeres n zu besserer Diskrimination
- Aber: mehr Null-Schätzungen \rightarrow führt zu Totalversagen
- Beispiel: Shakespeare als Korpus
 - $N=884,647$ tokens, $V=29,066$
 - 300,000 verschiedene Bigramme von $V^2 = 844m$ möglichen
 - 99.96% in Bigramtabelle sind Nullen
- Größeres Korpus kann besser funktionieren als kleineres
- Aber auch in einem großen Korpus werden wir nicht alles vorfinden (Chomsky, Zipfsches Gesetz)
- Mit Maximum Likelihood können wir nur “nacherzählen” \rightarrow Smoothing nötig

Smoothing

Reevaluation von Nullwahrscheinlichkeiten sowie n -grams mit niedrigen Wahrscheinlichkeiten, wobei ihnen höhere Werte zugewiesen werden.

Gegeben

$$P(w | \text{denied}, \text{the})$$

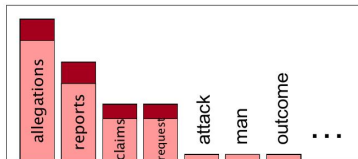
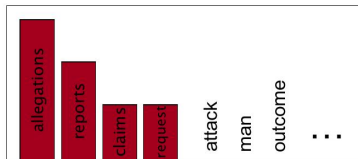
3 allegations

2 reports

1 claims

1 request

7 total



$$P_{+1}(w_i) = \frac{f(w_i) + 1}{N + V}$$

$$P_{+1}(w_2|w_1) = \frac{f(w_1, w_2) + 1}{f(w_1) + V}$$

$$P_{+1}(w_1, w_2 \dots w_n) = \frac{f(w_1, w_2, \dots w_n) + 1}{N + V^n}$$

- V ist Größe Vokabular(z.B., 20,000 Wörter);
- N ist Korpusgröße
- Bayesian Schätzung mit uniformer a priori Wahrscheinlichkeit

Sei X eine Zufallsvariable, die M Werte annehmen kann. Wir nennen den Wertebereich Ω_X . Ich schätze (durch N Experimente) die Wahrscheinlichkeiten der Werte. Sei $f(x)$ die Zählung für einen möglichen Wert x der Zufallsvariable. Laplace-Smoothing heißt dann:

$$P_{+1}(X = x) = \frac{f(x) + 1}{N + M}$$

Damit ist wieder eine Wahrscheinlichkeitsverteilung gegeben denn:

$$\sum_{x \in \Omega_X} P_{+1}(x) = \sum_{x \in \Omega_X} \left(\frac{f(x) + 1}{N + M} \right) = \frac{1}{N + M} \cdot \left(\sum_{x \in \Omega_X} f(x) + \sum_{x \in \Omega_X} 1 \right) = 1$$

Wir nehmen an, wir haben ein Vokabular von 10,000 Worten und einen Trainingskorpus von 1m Wörtern. Das Unigram *data* kommt zehnmal vor. Achtmal ist das nächste Wort *sparseness*.

$$P_{MLE}(sparseness|data) = ??$$

$$P_{+1}(sparseness|data) = ???$$

$$P_{+\lambda}(w_i) = \frac{f(w_i) + \lambda}{N + \lambda V}$$

$$P_{+\lambda}(w_2|w_1) = \frac{f(w_1, w_2) + \lambda}{f(w_1) + \lambda V}$$

$$P_{+\lambda}(w_1, w_2 \dots w_n) = \frac{f(w_1, w_2 \dots w_n) + \lambda}{N + \lambda V^n}$$

- $0 < \lambda < 1$
- Wie soll man λ schätzen?
- Immer noch vokabularabhängig
- Für n -gram Modelle ist Laplace/Lidstone-Smoothing ungeeignet
- Wenn man weniger Nullen hat, dann kann es benutzt werden
→ nächste Woche für Textklassifikation

Grundidee

Wenn wir die langen n -gramme nicht oder nicht oft genug sehen, dann vielleicht die kurzen

Interpolation: Kombiniere Unigram, Bigram und Trigram-Modelle.
Einfachste Version:

$$\begin{aligned}P_{intpol}(w_n|w_{n-2}w_{n-1}) &= \lambda_1 P(w_n|w_{n-2}w_{n-1}) \\ &+ \lambda_2 P(w_n|w_{n-1}) \\ &+ \lambda_3 P(w_n)\end{aligned}$$

mit $\sum_i \lambda_i = 1$ und $0 < \lambda_i < 1$ für $i = 1, 2, 3$

Teile Korpus in

- 1 Trainingskorpus. Für Originalschätzungen. Größter Teil des Korpus
- 2 Held-out Korpus/Development Set.

Fixiere n -gram Wahrscheinlichkeiten auf Trainingskorpus. Suche Lambdas, die die Wahrscheinlichkeiten auf Held-Out Korpus maximieren.

Maximiere (m Länge des held-out Korpus)

$$P_{interpol}(w_1, \dots, w_m) = \prod_{i=1}^m P_{interpol}(w_i | w_{i-2} w_{i-1})$$

Church und Gale (1991):

- Teilten 22m Korpus (AP newswire) in Training und held-out
- Für jedes Bigram im Trainingskorpus, schaue die wirkliche Häufigkeit im Held-out Korpus an
- Berechne Mittelwerte für Häufigkeitsklassen

Training	Held-Out
0	.0000270
1	0.448
2	1.25
3	2.24
4	3.23
5	4.21
6	5.23
7	6.21
8	7.21
9	8.26

$$P_{AD}(w_i|w_{i-1}) = \frac{\max(f(w_{i-1}w_i) - d, 0)}{f(w_{i-1})} + \lambda_{w_{i-1}}P(w_i)$$

- Eine Kombination von einfachem Discounting und Interpolation.
- d könnte im vorherigen Beispiel 0.75 sein
- Ist das Unigram am Ende der beste Backoff? Nein: Kneser-Ney (Optional)

Sollte $p(w_{i-1})$ einfach wie immer geschätzt werden?

- Vervollständige *I don't want to visit*
- Vielleicht ist *Francisco* häufiger als *Leavenworth*. Ist das trotzdem beste Ergänzung?

Benutze statt $p(w)$ besser $p_{Continuation}(w)$: Wie wahrscheinlich ist w als neue Fortsetzung? Nach wie vielen verschiedenen Wörtern (Types) kommt w vor?

$$p_{\text{continuation}}(\mathbf{w}) \propto |\mathbf{w}_{i-1} : f(\mathbf{w}_{i-1}, \mathbf{w}) > 0|$$

Muss normalisiert werden durch alle Bigramtypen:

$$p_{\text{continuation}}(\mathbf{w}) = \frac{|\mathbf{w}_{i-1} : f(\mathbf{w}_{i-1}, \mathbf{w}) > 0|}{|\{(\mathbf{w}_{j-1}, \mathbf{w}_j) : f(\mathbf{w}_{j-1}, \mathbf{w}_j) > 0\}|}$$

oder Normalisierung durch die Zahl aller Worte, die einem anderen vorangehen

$$p_{\text{continuation}}(\mathbf{w}) = \frac{|\mathbf{w}_{i-1} : f(\mathbf{w}_{i-1}, \mathbf{w}) > 0|}{\sum_{\mathbf{w}'} |f(\mathbf{w}'_{i-1}, \mathbf{w}') > 0|}$$

$$P_{KN}(w_i|w_{i-1}) = \frac{\max(f(w_{i-1}, w_i) - d, 0)}{f(w_{i-1})} + \lambda_{w_{i-1}} P_{\text{continuation}}(w_i)$$

λ ist eine Normalisierungskonstante, damit sich alles wieder auf 1 aufaddiert:

$$\lambda_{w_{i-1}} = \frac{d}{f(w_{i-1})} |w : f(w_{i-1}, w) > 0|$$

wobei der zweite Multiplikator: Anzahl der Worte, die auf w_{i-1} folgen können = Anzahl der Worttypen, für die wir einen Discount vorgenommen haben.

Wenn man Kneser-Ney auf höhere n-gramme anwendet, dann verkompliziert sich die Formel noch durch Rekursion (siehe J & M, Kapitel 3)

Alle bisherigen Methoden versuchen aus gegebenen Daten mehr herauszuholen. **Wie stehts mit mehr Daten?**

Idee: Häufigkeiten von $w_1 \dots w_n$ aus dem Web (Zhu :Rosenfeld, 2001)

$$\hat{P}_{web}(w_3|w_1, w_2) = \frac{f_{web}(w_1, w_2, w_3)}{f_{web}(w_1, w_2)}$$

- Hier nur für Nulltrigrams oder wenig gesehenen Trigrams
- Keller and Lapata 2004: Web-Häufigkeiten für **alle** n -gramme
- Wie bekommen wir Webhäufigkeiten?

Suchanfragen an Suchmaschine

- n -gram als Phrase an Suchmaschine
- wird genaue Phrasensuche ausführen
- Wird Anzahl der Webseiten ausgegeben \longrightarrow Schätzung

Was sind die Nachteile des Webs und/oder dieser Methode?

	Types	Nicht abgedeckt von			
		Alta Vista	Lycos	Google	Korpus
1g	327	0	0	0	8
2g	462	4	5	4	68
3g	453	46	46	23	189

Table 1: Novel n -gram types in 24 news sentences (Zhu and Rosenfeld, 2001)

- Korpus: 103 M newswire Korpus
- Testsätze noch nicht indiziert

- Markovmodelle schätzen Satzwahrscheinlichkeiten ab
- MLE Schätzung nicht brauchbar
- Einfaches Laplace-Smoothing vergibt zu viel Wahrscheinlichkeitsmasse an ungesehene Fälle
- Besser: Interpolation, Absolute Discounting, Kneser-Ney
- Im Moment beste Smoothingmethode: Kneser-Ney (Weiterentwicklung von Absolute Discounting mit Interpolation)
- * Jurafsky und Martin. Kapitel 3
- Übungsblatt 5

- 1 N-gram Modelle: Grundlagen
- 2 Schätzung von n -gram Modellen und Smoothing
- 3 Praxisprobleme und Intrinsische Evaluation**
- 4 Einfache Anwendungen: Confusion Sets und Sprachidentifikation

Oft wissen wir das gesamte Vokabular nicht im Voraus: Out of Vocabulary Words

- kreiere ein Token $\langle UNK \rangle$
- Training:
 - 1 Kreiere fixes Vokabular V
 - 2 Jedes Wort nicht in V wird durch $\langle UNK \rangle$ ersetzt
 - 3 Trainieren normal
- Beim testen: Benutze $\langle UNK \rangle$ Wahrscheinlichkeiten für jedes Wort, das wir nicht kennen

$$P(I \text{ drank white wine}) = \frac{3789}{1m} \cdot \frac{19}{1m} \cdot \frac{360}{1m} \cdot \frac{67}{1m}$$

Das werden aber kleine Zahlen \longrightarrow Rundungsprobleme
Mache alles im *log*-Raum

Berechne

$$\log P(w_1 \dots w_n) = \log p_1 \cdot p_2 \dots p_l = \log p_1 + \log p_2 + \dots + \log p_l$$

Intuition

Ein Textmodell ist besser, wenn es eine höhere Wahrscheinlichkeit den Wörtern und Wortsequenzen zuordnet, die im Testset wirklich vorkommen. D.H. Modell A ist besser als B , wenn es das Testkorpus als wahrscheinlicher ansieht. Man ist weniger überrascht.

Methode:

- 1 Spalte Korpus in Trainings- und Testkorpus. (Warum?)
- 2 Schätze Parameter des Modells auf Trainingskorpus. Dieser ist oft noch mal aufgespalten in Training und Held-out.
- 3 Teste Performanz auf Testkorpus. Messe zum Beispiel (Kreuz)Entropie und Perplexität des Modells auf Testkorpus.

Cross-entropy und Perplexität

Cross-entropy zwischen Sprache L und Modell p der Sprache gemessen auf Korpus $W := w_1, \dots, w_N$

$$H(W, p) \approx -\frac{1}{N} \log p(w_1, \dots, w_N)$$

Perplexität:

$$PP(W) := 2^{H(W, p)} = 2^{-\frac{1}{N} \log p(w_1, \dots, w_N)} = \sqrt[N]{\frac{1}{p(w_1 \dots w_N)}}$$

Kettenregel:

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{p(w_i | w_1 \dots w_{i-1})}}$$

Bigrams:

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{p(w_i | w_{i-1})}}$$

- Nehmen wir an ein Satz bestehe aus aus zufälliger Anordnung von Ziffern Null bis 9.
- Perplexität eines Unigrammodells, das jede Ziffer als gleichwahrscheinlich annimmt $p(d) = 1/10$ für alle d

$$\begin{aligned} PP(W) &= p(w_1, \dots, w_N)^{-\frac{1}{N}} \\ &= \left(\frac{1}{10}\right)^{-\frac{1}{N}} \\ &= 10 \end{aligned}$$

Bei zufälligen gleichverteilten Variablen kann man nicht besser werden. Sonst aber schon!!!

Korpus:

- $\langle s \rangle$ there is a big house $\langle /s \rangle$
- $\langle s \rangle$ i buy a house $\langle /s \rangle$
- $\langle s \rangle$ they buy the new house $\langle /s \rangle$

Testsatz: $\langle s \rangle$ *they buy a big house* $\langle /s \rangle$

Unter Bigrammodell (kein smoothing):

$$p(S) = 3/20 \cdot 0.333 \cdot 1 \cdot 0.5 \cdot 0.5 \cdot 1 \cdot 1 = 0.01245$$

Kreuzentropie:

$$-\frac{1}{7} \log p(S) = 0.891$$

Perplexität: $PP(S) = 1.8544$

Training 38 m Wörter, Testset 1.5 Millionen Wörter, WSJ

Perplexität:

Unigram	Bigram	Trigram
962	170	109

Tools und Daten:

- **SRILM**: <http://www.speech.sri.com/projects/srilm/>
- **KenLM**: <https://kheafield.com/code/kenlm/>
- Natürlich auch Implementationen in NLTK
- Google n -gram Datensatz: Korpus über 1 Billion Wörter, Zahlen für alle n -grams von 1 bis 5, die mindestens 5mal vorkommen.

serve as the incoming	92
serve as the incubator	99
serve as the independent	794
serve as the index	223
serve as the indication	72

- 1 N-gram Modelle: Grundlagen
- 2 Schätzung von n -gram Modellen und Smoothing
- 3 Praxisprobleme und Intrinsische Evaluation
- 4 Einfache Anwendungen: Confusion Sets und Sprachidentifikation**

Wortwahl, wenn es eine kleine Anzahl an Möglichkeiten gibt

Beispiele:

- Handschriftenerkennung: *a big brown b....r*
- MT:
 - ① Ich glaube diese Geschichte nicht.
 - ② I don't believe this { *history,story,tale,saga,strip.* }
- Rechtschreibkorrektur
- Ordnung in Generierung

Anwendungen sind eine gute Möglichkeit, die Güte eines Modells abzuschätzen

Resultat	Nichtwörter	Wörter
Grund	typographisch	kognitiv
Fehleranzahl	einer	mehrere

Minimum Edit Distance war gut für das Vorschlagen von Korrekturen bei kleiner Fehleranzahl. Nicht gut für das Auswählen der richtigen Korrektur.

- **non-word error detection**: Erkenne Nichtwörter
- **isolated-word error correction**: korrigiere Nichtwörter ohne Kontext (Min Edit, isolierte Häufigkeit der Korrektur)
- **context-dependent error detection and correction**: entdecke und korrigiere real word errors.

Confusion sets

Wortpaare oder Wortmengen, die leicht füreinander falsch geschrieben werden

principle	principal	
then	than	
affect	effect	
quite	quiet	
peace	piece	
cite	site	sight
you're	your	
their	there	they're
weather	whether	

*I don't know **weather** I want to go to the party.* → richtig oder falsch?

$w_{-n} w_{-n+1} \dots w_{-1} t w_1 w_2 \dots w_m$ → t richtig oder falsch?

äquivalent zu

I don't know ? I want to go to the party. → wähle *weather* oder *whether*

$w_{-n} w_{-n+1} \dots w_{-1} ? w_1 w_2 \dots w_m$ → welches Mitglied t des confusion set soll die Lücke ausfüllen

Geg. $w_{-n} w_{-n+1} \dots w_{-1} ? w_1 w_2 \dots w_m$ und $t_1 \dots t_j$ als confusion set

- Unigram model: $\max_{t_i \in t_1 \dots t_j} f(t_i)$
- Linkes-bigram: $\max_{t_i \in t_1 \dots t_j} f(w_1, t_i)$
- Rechtes bigram : $\max_{t_i \in t_1 \dots t_j} f(t_i, w_1)$
- Mittleres trigram : $\max_{t_i \in t_1 \dots t_j} f(w_{-1}, t_i, w_1)$
- Linkes trigram: $\max_{t_i \in t_1 \dots t_j} f(w_{-2}, w_{-1}, t_i)$
- Rechtes trigram: $\max_{t_i \in t_1 \dots t_j} f(t_i, w_1, w_2)$

Modelle geschätzt auf Korpus C

Testsatz: *I don't know ? I want to go to the party*; confusion set
whether, wheather

- Unigram: $f(\textit{whether})$ vs. $f(\textit{weather})$
- Linkes Bigram: $f(\textit{know whether})$ vs. $f(\textit{know weather})$
- Rechtes bigram: $f(\textit{whether I})$ vs $f(\textit{weather I})$
- Mittleres Trigram ?
- Rechtes trigram ?
- Linkes trigram?

Geg. confusion set S mit Elementen $t_1 \dots t_j$

- 1 Nimm Korpus $C1$ and extrahiere alle Kookkurenzen $t_1 \dots t_j$ mit Kontext
- 2 Ersetze durch Lücke und versuche Original vorherzusagen
- 3 Messe Akkuratheit auf diese automatisch generiertem Testset

$C1$ muss von C getrennt sein

Aus Korpus *C1* und confusion set { *whether, wheather* } leite Testset ab:

Kontext	Original	Vorhersage?
I don't know ? I want to	whether	{ whether, weather }
The ? is pretty bad today	weather	{ whether, weather }
? he likes me or not I can't say	whether	{ whether, weather }
I like sunny ?	weather	{ whether, weather }

Prädiktion via n-gram model auf Korpus *C*

18 confusion sets, 2056 Beispiele, Schätzungen BNC; Keller und Lapata (2005)

Model	BNC
Unigram	71.78
Linkes bigram	83.29
Rechtes bigram	81.70
Mittleres trigram	78.14
Linkes trigram	74.72
Rechtes trigram	73.46

18 confusion sets, 2056 Beispiele

Model	Web	BNC
Unigram	74.79	71.78
Linkes bigram	86.13	83.29
Rechtes bigram	82.36	81.70
Mittleres trigram	89.35	78.14
Linkes trigram	89.24	74.72
Rechtes trigram	83.71	73.46

- 1 Unüberwachte Methode mit sehr kurzem Kontext: 89.35%
- 2 Web viel besser als BNC, trotz vieler theoretischer Nachteile
- 3 Gute überwachte Modelle erreichen 94% und mehr %
- 4 Methode kann auch auf confusion sets, die nichts mit Rechtschreibung zu tun haben, generalisiert werden (*until/by*)

N-grams Anwendungen

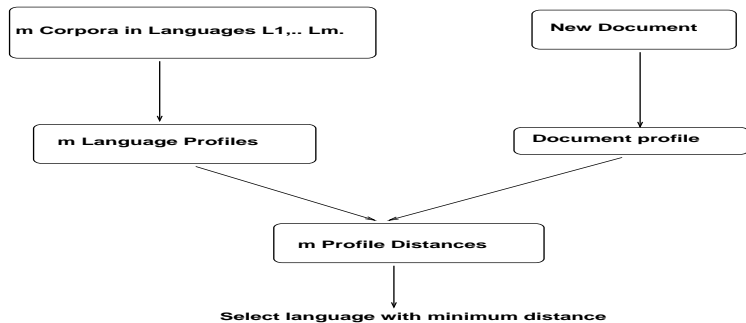
- Oft gut für Wortwahl zwischen kleiner Menge an Möglichkeiten
- Braucht smoothing und/oder sehr große Korpora
- Kann gute, billige (unüberwachte) Alternative zu überwachtem maschinellem Lernen darstellen

Lapata and Keller (2005): Web-based models for Natural Language Processing. IN ACM Transactions on Speech and Language Processing

Aufgabe: gegeben eine Zeichenkette, welche (menschliche) Sprache ist es?

Vorschläge

Buchstaben n -gramme für Sprachidentifikation



- **Schritt 1:** Trainingsdaten für mehrere Sprachen
- **Schritt 2:** Häufigkeiten von n -grammen für Buchstaben n -grams aus Trainingsset
- Carver und Trenkle (1994): 1-5 gramme

Rang	n-gram
1	e
2	t
...	...
...	th
...	er
...	on
...	le
...	ing
...	and



Für Testdokument:

- **Schritt 3:** Profil für das Testdokument
- ▶ **Schritt 4:** Vergleiche Sprachprofile mit Testdocumentprofil
- **Schritt 5:** Wähle die Sprache mit dem ähnlichstem Profil

Rang	Sprachprofil	Testprofil
1	e	e
2	t	t
3	th	th
4	er	ing
5	on	on
6	le	er
7	ing	and
8	and	ed
...

Rang	Sprachprofil	Testprofil	Unterschied
1	e	e	0
2	t	t	0
3	th	th	0
4	er	ing	3
5	on	on	0
6	le	er	2
7	ing	and	1

3200 newsgroup Artikel in 8 Sprachen; leave-one-out classification

L	<300	<300	<300	> 300	> 300	> 300
Profil	100	200	300	100	200	300
Acc	92.9%	97.6%	98.6%	97.2%	99.5%	99.8%

L ist Artikellänge in Byte

- Funktioniert schon gut auf kurzen Artikeln
- fast perfekt für lange Artikel
- auch gut für verrauschte Daten
- Demo: <https://tools.wmflabs.org/textcatdemo/>

- * Jurafsky und Martin, Kapitel 3
- Frank Keller and Mirella Lapata: Using the Web to Obtain Frequencies for Unseen Bigrams, Computational Linguistics, 2003
- Zhu and Rosenfeld, 2001: Improving Trigram Language Modelling with the World Wide Web. In: *Proceedings of International Conference on Acoustics, Speech and Signal Processing*
- Cavnar and Trenkle (1994): N-gram-based Text Categorisation. Symposium on Document Analysis and Information Retrieval.
- Lapata and Keller (2005): Web-based models for Natural Language Processing. IN ACM Transactions on Speech and Language Processing