

Oberseminar

25.07.2013

Introduction to Kneser-Ney Smoothing on Top of Generalized Language Models for Next Word Prediction

Martin Körner

- Introduction
- Language Models
- Generalized Language Models
- Smoothing
- Progress
- Summary

- **Introduction**
- Language Models
- Generalized Language Models
- Smoothing
- Progress
- Summary

- Next word prediction: What is the next word a user will type?
- Use cases for next word prediction:
 - ◆ Augmentative and Alternative Communication (AAC)
 - ◆ Small keyboards (Smartphones)



- How do we predict words?
 1. Rationalist approach
 - Manually encoding information about language
 - “Toy” problems only
 2. Empiricist approach
 - Statistical, pattern recognition, and machine learning methods applied on corpora
 - Result: Language models

- Introduction
- **Language Models**
- Generalized Language Models
- Smoothing
- Progress
- Summary

- Language model: How likely is a sentence s ?
→ Probability distribution: $P(s)$
- Calculate $P(s)$ by multiplying conditional probabilities
- Example:

$P(\text{If you're going to San Francisco , be sure ...})$

=

$P(\text{you're} \mid \text{If}) * P(\text{going} \mid \text{If you're}) *$

$P(\text{to} \mid \text{If you're going}) * P(\text{San} \mid \text{If you're going to}) *$

$P(\text{Francisco} \mid \text{If you're going to San}) * \dots$

→ Empirical approach would fail

- Markov assumption [JM80]:
 - ♦ Only the last $n-1$ words are relevant for a prediction
 - ♦ Example with $n=5$:

$$P(\text{sure} \mid \text{If you're going to San Francisco , be}) \\ \approx P(\text{sure} \mid \text{San Francisco , be})$$

 Counts as a word

- n -gram: Sequence of length n with a count
 - ♦ E.g.: 5-gram:

If you're going to San 4

- Sequence naming:

$$w_1^{i-1} := w_1 w_2 \dots w_{i-1}$$

- Markov assumption formalized:

$$P(w_i | w_1^{i-1}) \approx P(w_i | \underbrace{w_{i-n+1}^{i-1}}_{n-1 \text{ words}})$$

- Instead of $P(s)$:

- ♦ Only one conditional probability $P(w_i | w_{i-n+1}^{i-1})$

- Simplify $P(w_i | w_{i-n+1}^{i-1})$ to $P(w_n | w_1^{n-1})$

$n-1$ words

$n-1$ words

Conditional probability with Markov assumption

$$\text{NWP}(w_1^{n-1}) = \arg \max_{w_n \in W} \left(P(w_n | w_1^{n-1}) \right)$$

Set of all words in the corpus

→ How to calculate the probability $P(w_n | w_1^{n-1})$?

- The easiest way:
 - ♦ Maximum likelihood:

$$P_{\text{ML}}(w_n | w_1^{n-1}) = \frac{c(w_1^n)}{c(w_1^{n-1})}$$

- ♦ Example:

$$P(\text{San} | \text{If you're going to}) = \frac{c(\text{If you're going to San})}{c(\text{If you're going to})}$$

- Introduction
- Language Models
- **Generalized Language Models**
- Smoothing
- Progress
- Summary

- Main idea:

- ◆ Insert wildcard words (*) into sequences

- Example:

- ◆ Instead of $P(\text{San} \mid \text{If you're going to})$:

- $P(\text{San} \mid \text{If } * * *)$

- $P(\text{San} \mid \text{If } * * \text{ to})$

- $P(\text{San} \mid \text{If } * \text{ going } *)$

- $P(\text{San} \mid \text{If } * \text{ going to})$

- $P(\text{San} \mid \text{If you're } * *)$

- ...

Length: 5, Wildcard words: 2

→ Aggregate results

- ◆ Separate different types of GLMs based on:

1. Sequence length
2. Number of wildcard words

- Data sparsity of n -grams
 - ◆ “If you’re going to San” is seen less often than for example
“If * * to San”

→ Question: Does that really improve the prediction?

- ◆ Result of evaluation: Yes

... but we should use smoothing for language models

- Introduction
- Language Models
- Generalized Language Models
- **Smoothing**
- Progress
- Summary

- Problem: Unseen sequences
 - Try to estimate probabilities of unseen sequences
 - ◆ Probabilities of seen sequences need to be reduced

- Two approaches:
 1. Backoff smoothing
 2. Interpolation smoothing

- If sequence unseen: use shorter sequence
 - ♦ E.g.: if $P(\text{San} \mid \text{going to}) = 0$ use $P(\text{San} \mid \text{to})$

Higher order
probability

$$P_{back}(w_n | w_i^{n-1}) = \begin{cases} \tau(w_n | w_i^{n-1}) & \text{if } c(w_i^n) > 0 \\ \gamma * P_{back}(w_n | w_{i+1}^{n-1}) & \text{if } c(w_i^n) = 0 \end{cases}$$

Weight

Lower order
probability (recursive)

- Always use shorter sequence for calculation

$$P_{inter}(w_n | w_i^{n-1}) = \underbrace{\tau(w_n | w_i^{n-1})}_{\text{Higher order probability}} + \underbrace{\gamma}_{\text{Weight}} * \underbrace{P_{inter}(w_n | w_{i+1}^{n-1})}_{\text{Lower order probability (recursive)}}$$

- Seems to work better than backoff smoothing

- Interpolated smoothing
- Idea: Improve lower order calculation
- Example: Word visiting unseen in corpus

$$P(\text{Francisco} \mid \text{visiting}) = 0$$

→ Normal interpolation: $0 + \gamma * P(\text{Francisco})$

$$P(\text{San} \mid \text{visiting}) = 0$$

→ Normal interpolation: $0 + \gamma * P(\text{San})$

Result: Francisco is as likely as San at that position

Is that correct?

→ Difference between Francisco and San?

Answer: Number of different contexts

- For lower order calculation:

- ♦ Don't use $c(w_n)$
- ♦ Instead: Number of different bigrams the word completes:

$$N_{1+}(\bullet w_n) := |\{w_{n-1} : c(w_{n-1}^n) > 0\}|$$

Count

- ♦ Or in general:

$$N_{1+}(\bullet w_{i+1}^n) = |\{w_i : c(w_i^n) > 0\}|$$

- In addition:

- ♦ $N_{1+}(\bullet w_{i+1}^{n-1} \bullet) = \sum_{w_n} N_{1+}(\bullet w_{i+1}^n)$
- ♦ $N_{1+}(w_i^{n-1} \bullet) = |\{w_n : c(w_i^n) > 0\}|$

- ◆ Highest order calculation:

$$P_{\text{KN}}(w_n | w_i^{n-1}) = \frac{\max\{c(w_i^n) - D, 0\}}{c(w_i^{n-1})} + \underbrace{\frac{D}{c(w_i^{n-1})} N_{1+}(w_i^{n-1})}_{\text{Lower order weight}} \cdot \underbrace{P_{\text{KN}}(w_n | w_{i+1}^{n-1})}_{\text{Lower order probability (recursion)}}$$

Assure positive value → $\max\{c(w_i^n) - D, 0\}$
 count → $c(w_i^n)$
 Discount value $0 \leq D \leq 1$ → $- D$
 Total counts → $c(w_i^{n-1})$

- ◆ Lower order calculation:

$$P_{\text{KN}}(w_n | w_i^{n-1}) = \frac{\max\{N_{1+}(\bullet w_i^n) - D, 0\}}{N_{1+}(\bullet w_i^{n-1} \bullet)} +$$

Assure positive value
Continuation count
Discount value

$$\frac{D}{N_{1+}(\bullet w_i^{n-1} \bullet)} N_{1+}(w_i^{n-1} \bullet) P_{\text{KN}}(w_n | w_{i+1}^{n-1})$$

Lower order weight
Lower order probability (recursion)

- ◆ Lowest order calculation: $P_{\text{KN}}(w_n) = \frac{N_{1+}(\bullet w_i^n)}{N_{1+}(\bullet w_i^{n-1} \bullet)}$

- Different discount values for different absolute counts
- Lower order calculation:

$$P_{\text{KN}}(w_n | w_i^{n-1}) = \frac{\max\{N_{1+}(\bullet w_i^n) - D(c(w_i^n)), 0\}}{N_{1+}(\bullet w_i^{n-1} \bullet)} +$$

$$\frac{D_1 N_1(w_i^{n-1} \bullet) + D_2 N_2(w_i^{n-1} \bullet) + D_{3+} N_{3+}(w_i^{n-1} \bullet)}{N_{1+}(\bullet w_i^{n-1} \bullet)} P_{\text{KN}}(w_n | w_{i+1}^{n-1})$$

- State of the art (since 15 years!)

- We can use all smoothing techniques on GLMs as well!

- Small modification:

E.g: $P(\text{San} \mid \text{If } * \text{ going } *)$

Lower order sequence :

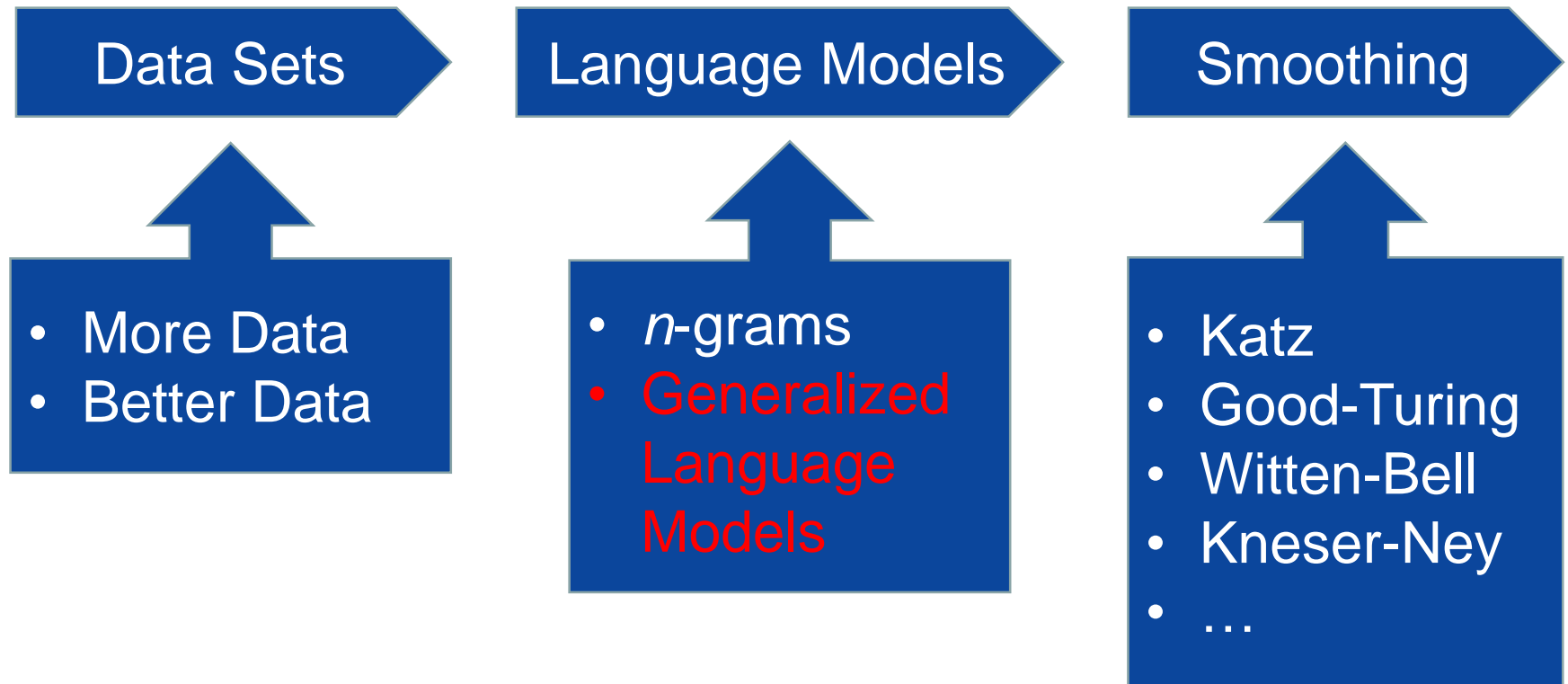
- Normally: $P(\text{San} \mid * \text{ going } *)$
- Instead use $P(\text{San} \mid \text{going } *)$

- Introduction
- Language Models
- Generalized Language Models
- Smoothing
- **Progress**
- Summary

- Done Yet:
 - ◆ Extract text from XML files
 - ◆ Building GLMs
 - ◆ Kneser-Ney and modified Kneser-Ney smoothing
 - ◆ Indexing with MySQL

- ToDo's
 - ◆ Finish evaluation program
 - ◆ Run evaluation
 - ◆ Analyze results

- Introduction
- Language Models
- Generalized Language Models
- Smoothing
- Progress
- **Summary**



Thank you for your attention!

Questions?

■ Images:

- ◆ Wheelchair Joystick (Slide 4):
http://i01.i.aliimg.com/img/pb/741/422/527/527422741_355.jpg
- ◆ Smartphone Keyboard (Slide 4):
https://activecaptain.com/articles/mobilePhones/iPhone/iPhone_Keyboard.jpg

■ References:

- ◆ **[CG98]**: Stanley Chen and Joshua Goodman. An empirical study of smoothing techniques for language modeling. Technical report, Technical Report TR-10-98, Harvard University, August, 1998.
- ◆ **[JM80]**: F. Jelinek and R.L. Mercer. Interpolated estimation of markov source parameters from sparse data. In Proceedings of the Workshop on Pattern Recognition in Practice, pages 381–397, 1980.
- ◆ **[KN95]**: Reinhard Kneser and Hermann Ney. Improved backing-off for m-gram language modeling. In Acoustics, Speech, and Signal Processing, 1995. ICASSP-95., 1995 International Conference on, volume 1, pages 181–184. IEEE, 1995.