

Colloquium

12.09.2013

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

Martin Körner

- Introduction
- Language Models
- Generalized Language Models
- Smoothing
- Implementation
- Future Work
- Summary

- **Introduction**
- Language Models
- Generalized Language Models
- Smoothing
- Implementation
- Future Work
- 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
- Implementation
- Future Work
- Summary

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

- ♦ Example:

$$\begin{aligned} P(\text{If you're going to San Francisco , be sure ...}) = \\ P(\text{you're} \mid \text{If}) * P(\text{going} \mid \text{If you're}) * \\ P(\text{to} \mid \text{If you're going}) * P(\text{San} \mid \text{If you're going to}) * \\ P(\text{Francisco} \mid \text{If you're going to San}) * \dots \end{aligned}$$

- Instead of $P(s)$ only one conditional probability:

$$P(\text{Next Word} \mid \text{Current Sentence})$$

→ Empiricist approach fails

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

$$P(\text{Next Word} \mid \text{Last } n - 1 \text{ Words})$$

- ♦ Example with $n=5$:

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

 Counts as a word

- n -gram: Word 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:

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

Current sentence

$n-1$ words

$$P(w_n | w_1^{n-1})$$


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

- 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:

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

- ♦ Still: Problem with data sparsity

- Introduction
- Language Models
- **Generalized Language Models**
- Smoothing
- Implementation
- Future Work
- 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**
- Implementation
- Future Work
- 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

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

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

→ Normal backoff: $\gamma * P(\text{Francisco})$

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

→ Normal backoff: $\gamma * P(\text{San})$

Also: San and Francisco have the same counts

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:

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

- Modified Kneser-Ney [CG98]: Different calculation of higher order discounts

- 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
- **Implementation**
- Future Work
- Summary

- SRI Language Modeling Toolkit, German Wiki (2.6 GB text)

Total: 30 Seconds

```
 1 [|||||||||||||||||||] 43.7%] Tasks: 98, 169 thr; 2 running
 2 [|||||||||||||||||||] 64.2%] Load average: 0.86 1.05 1.08
 3 [|||||] 8.7%] Uptime: 00:52:44
 4 [|||||] 9.2%]
Mem[|||||||||||||||||] 4181/15750MB]
Swp[||] 224/8191MB]

PID USER      PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
10383 martin    20   0 3798M 3778M 1248  R 100. 24.0  0:29.98 ./ngram-count -or
```

Process: 30 Seconds

- SRI Language Modeling Toolkit, German Wiki (2.6 GB text)

Total: 60 Seconds

```
 1  [|| 2.0%] Tasks: 98, 168 thr; 3 running
 2  [||||| 12.8%] Load average: 0.92 1.04 1.07
 3  [||||||| 21.1%] Uptime: 00:53:14
 4  [||||||||| 66.2%]
Mem [||||||||| 6776/15750MB]
Swp [|| 224/8191MB]

  PID USER   PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
10383 martin  20   0 6387M 6367M 1248  R 100. 40.4  1:00.03 ./ngram-count -or
```

Process: 60 Seconds

- SRI Language Modeling Toolkit, German Wiki (2.6 GB text)

Total: 120 Seconds

```
 1  [||] 2.0%] Tasks: 98, 192 thr; 2 running
 2  [|||||||||||||||||||||] 52.3%] Load average: 0.97 1.04 1.07
 3  [|||||||||||||||||||||] 48.7%] Uptime: 00:54:15
 4  [|] 0.7%]
Mem [|||||||||||||||||||||] 11245/15750MB]
Swp [||] 224/8191MB]

  PID USER   PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
10383 martin  20   0 10.5G 10.5G 1248  R 100. 68.2  2:00.18 ./ngram-count -or
```

Process: 120 Seconds

- SRI Language Modeling Toolkit, German Wiki (2.6 GB text)

Total: 180 Seconds

```
1  [|||] 6.0%] Tasks: 99, 195 thr; 2 running
2  [||] 1.4%] Load average: 1.22 1.10 1.09
3  [|||||||||||||||||||||||||||||||||]100.0%] Uptime: 00:55:15
4  [||] 3.3%]
Mem[|||||||||||||||||||||||||||||]15347/15750MB]
Swp[||] 224/8191MB]

  PID USER      PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
10383 martin    20   0 14.5G 14.5G 1232  R 100. 94.3  3:00.43 ./ngram-count -or
```

Process: 180 Seconds

- SRI Language Modeling Toolkit, German Wiki (2.6 GB text)

Total: 1560 Seconds

```
 1  [|||]          4.6%]      Tasks: 99, 196 thr; 1 running
 2  [|]           1.3%]      Load average: 1.27 1.28 1.21
 3  [|||]        5.4%]      Uptime: 01:18:02
 4  [||]         4.0%]
Mem[|||||||||||||||||||||||||15417/15750MB]
Swp[|||||]       1918/8191MB]

  PID USER      PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
10383 martin    20   0 15.7G 14.0G 364  D  3.3  91.3  4:03.31 ./ngram-count -or
```

Process: 240 Seconds

- SRI Language Modeling Toolkit, German Wiki (2.6 GB text)

Total: 6540 Seconds

```
1  [||] 2.0%] Tasks: 113, 255 thr; 1 running
2  [||] 2.0%] Load average: 1.39 1.37 1.30
3  [||] 4.0%] Uptime: 02:31:13
4  [|||] 5.3%]
Mem[|||||||||||||||||||||15236/15750MB]
Swp[|||||||||||||||||4111/8191MB]

PID USER      PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
10383 martin    20   0 16.8G 13.1G 364  D  0.0  85.3  6:00.51 ./ngram-count -or
```

Process: 360 Seconds

- English Wikipedia: 6.6 GB

- Kyoto Language Modeling Toolkit, German Wiki (2.6 GB)

```
martin@martin-ThinkPad:~/glm/kylm$ java -cp kylm-0.0.7.jar kylm.main.CountNgrams de-training.txt de-ml5.arpa -n 5 -ml
```

Guess what is going to happen
(Hint: It's written in Java)

- Kyoto Language Modeling Toolkit, German Wiki (2.6 GB)

```
martin@martin-ThinkPad:~/glm/kylm$ java -cp kylm-0.0.7.jar kylm.main.CountNgrams de-training.txt de-ml5.arpa -n 5 -ml
Exception in thread "main" java.lang.OutOfMemoryError: Java heap space
    at java.util.Arrays.copyOf(Arrays.java:2219)
    at java.util.Vector.grow(Vector.java:262)
    at java.util.Vector.ensureCapacityHelper(Vector.java:242)
    at java.util.Vector.insertElementAt(Vector.java:597)
    at kylm.model.ngram.BranchNode.getChild(Unknown Source)
    at kylm.model.ngram.NgramLM.countNgrams(Unknown Source)
    at kylm.model.ngram.NgramLM.trainModel(Unknown Source)
    at kylm.main.CountNgrams.main(Unknown Source)
martin@martin-ThinkPad:~/glm/kylm$
```

java.lang.OutOfMemoryError: Java heap space



- We need ...
 - ... to work on large corpora (6 GB +)
 - ... to work on many different corpora
 - ... Generalized Language Models
 - ... Modified Kneser-Ney Smoothing

- How we build n -grams and generalized n -grams:
 1. Store different sequences in many (small) files
 2. Sort the files
 3. Count the sequences (easy, because they're sorted)
 4. Store counts
- Some numbers for English Wikipedia:
 - ◆ Size of text file: 6.6 GB

| Type | Different sequences | Size |
|--------------------------|---------------------|--------|
| 5-grams | 862.489.285 | 25 GB |
| (1..5)-grams | 2.004.173.723 | 51 GB |
| Generalized (1..5)-grams | 10.057.876.269 | 217 GB |

1. Calculate different counts: $N_{1+}(\bullet w_{i+1}^n)$, $N_{1+}(w_{i+1}^{n-1} \bullet)$...
2. Calculate lowest order probability
3. Calculate second-lowest order probability
4. Aggregate lowest and second-lowest order probabilities

→ Problem: Sorting of those sequences

Example: Higher order: to San

Lower order: San

Solution: Sort higher order seq. starting at second word

5. Repeat steps 3 and 4 with higher order sequences

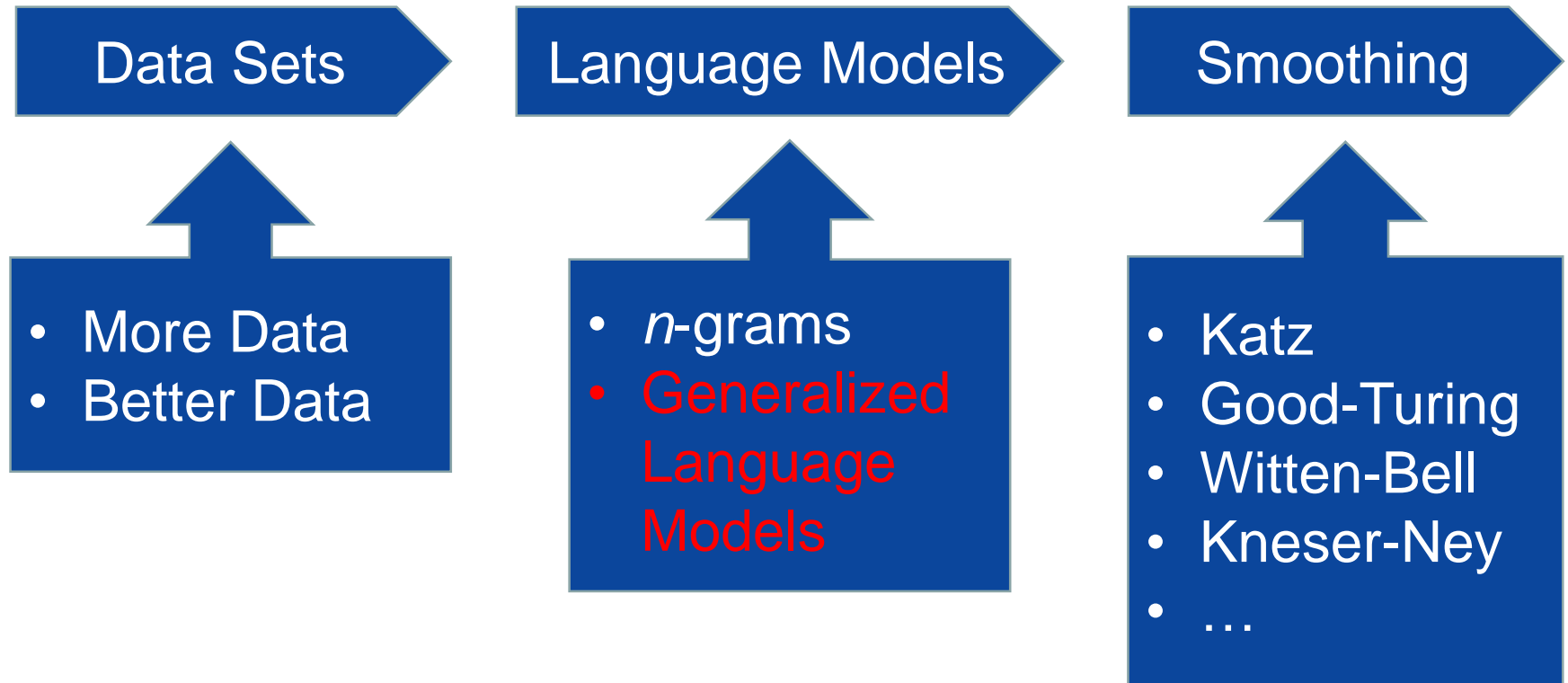
- Introduction
- Language Models
- Generalized Language Models
- Smoothing
- Implementation
- **Future Work**
- Summary

- Evaluation
 - ◆ Evaluating next word prediction
 - ◆ More general metrics like perplexity and cross-entropy

- Different way of backing off at Kneser-Ney for Generalized Language Models (idea from Thomas Gottron)

- Code cleanup
- Better documentation
- ...

- Introduction
- Language Models
- Generalized Language Models
- Smoothing
- Implementation
- Future Work
- **Summary**



Find my bachelor thesis online at:

<http://mkrnr.de/bt>


That's MKRNR

Also, the slides from ...

... last Oberseminar talk:

<http://mkrnr.de/c>

... today:

<http://mkrnr.de/d>

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.