

# Natural Language Processing: Foundations

Section 3 — n-Gram Language Models (Motivation)

# Language Models — Motivation

- Which sentence makes more sense?  $S_1$  or  $S_2$ ?

$S_1$ : "on guys all I of noticed sidewalk three a sudden standing the"

Example 1:

$S_2$ : "all of a sudden I noticed three guys standing on the sidewalk"

$S_1$ : "the role was played by an **aeressacross** famous for her comedic timing"

Example 2:

$S_2$ : "the role was played by an **aeressactress** famous for her comedic timing"

- But why?
  - Probability of  $S_2$  higher than of  $S_1$ :  $P(S_2) > P(S_1)$

→ **Language Models** — Assigning probabilities to a sentence, phrase (or word)

# Language Models — Basic Idea

- 2 basic notions of probabilities

(1) **Probability of a sequence of words**  $W$

$$P(W) = P(w_1, w_2, w_3, \dots, w_n)$$

Example:  $P(\text{"remember to submit your assignment"})$

$\begin{matrix} ^1 \\ | \\ \backslash \\ \text{r} \end{matrix} \quad \begin{matrix} ^2 \\ | \\ \backslash \\ \text{e} \\ \text{m} \\ \text{e} \\ \text{m} \\ \text{e} \\ \text{r} \\ \text{y} \\ \text{t} \\ \text{t} \\ \text{u} \\ \text{b} \\ \text{i} \\ \text{s} \\ \text{u} \\ \text{b} \\ \text{m} \\ \text{e} \\ \text{n} \\ \text{c} \\ \text{h} \\ \text{o} \\ \text{r} \\ \text{g} \end{matrix} \quad \begin{matrix} ^3 \\ | \\ \backslash \\ \text{a} \\ \text{s} \\ \text{s} \\ \text{e} \\ \text{c} \\ \text{h} \\ \text{a} \\ \text{g} \\ \text{m} \\ \text{e} \\ \text{n} \\ \text{t} \\ \text{t} \\ \text{u} \\ \text{r} \\ \text{e} \\ \text{m} \\ \text{e} \\ \text{m} \\ \text{e} \\ \text{r} \\ \text{y} \\ \text{t} \\ \text{t} \\ \text{u} \\ \text{b} \\ \text{i} \\ \text{s} \\ \text{u} \\ \text{b} \\ \text{m} \\ \text{e} \\ \text{n} \\ \text{c} \\ \text{h} \\ \text{o} \\ \text{r} \\ \text{g} \end{matrix} \quad \dots \quad \begin{matrix} ^n \\ | \\ \backslash \\ \text{a} \\ \text{s} \\ \text{s} \\ \text{e} \\ \text{c} \\ \text{h} \\ \text{a} \\ \text{g} \\ \text{m} \\ \text{e} \\ \text{n} \\ \text{t} \\ \text{t} \\ \text{u} \\ \text{r} \\ \text{e} \\ \text{m} \\ \text{e} \\ \text{m} \\ \text{e} \\ \text{r} \\ \text{y} \\ \text{t} \\ \text{t} \\ \text{u} \\ \text{b} \\ \text{i} \\ \text{s} \\ \text{u} \\ \text{b} \\ \text{m} \\ \text{e} \\ \text{n} \\ \text{c} \\ \text{h} \\ \text{o} \\ \text{r} \\ \text{g} \end{matrix}$

(2) **Probability of an upcoming word**  $w_n$

prob F-Y

$$P(w_n \mid w_1, w_2, w_3, \dots, w_{n-1})$$

Example:  $P(\text{"assignment"} \mid \text{"remember to submit your"})$

$P(X \mid Y)$   
↑  
 $Y$  already happened

In this section: How to calculate these probabilities?

# Language Models — Applications

- Language Models are fundamental for many NLP tasks
  - **Speech Recognition**  $P(\text{"we built this city on rock and roll"}) > P(\text{"we built this city on sausage rolls"})$
  - **Spelling correction**  $P(\text{"... has no mistakes"}) > P(\text{"... has no \underline{mistaek}}")$
  - **Grammar correction**  $P(\text{"... has improved"}) > P(\text{"... has \underline{improve}}")$
  - **Machine Translation**  $P(\text{"I went home"}) > P(\text{"I went to home"})$

# **Natural Language Processing: Foundations**

Section 3 — n-Gram Language Models (Sentence Probabilities)

# Probabilities of Sentences

(more generally: sequence of words)

$$\left. \begin{array}{l} P(\text{"remember to submit your assignment"}) \\ P(\text{"assignment"} | \text{"remember to submit your"}) \end{array} \right\} \rightarrow \text{How to calculate those probabilities?}$$

- Quick review: Chain Rule (allows the iterative calculation of joint probabilities)

- Chain rule for 2 random events:  $P(A_1, A_2) = P(A_2|A_1) \cdot P(A_1) = P(A_1|A_2) \cdot P(A_2)$
- Chain rule for 3 random events:  $P(A_1, A_2, A_3) = P(A_3|A_1, A_2) \cdot P(A_1, A_2)$   
 $= P(A_3|A_1, A_2) \cdot P(A_2|A_1) \cdot P(A_1)$
- ...

# Probabilities of Sentences

Unordered  
order arbitrary.

- Chain rule — generalization to  $N$  random events

$$\begin{aligned} P(A_1, \dots, A_N) &= P(A_1) \cdot P(A_2|A_1) \cdot P(A_3|A_{1:2}) \cdot \dots \cdot P(A_N|A_{1:N-1}) \\ &= \prod_{i=1}^N P(A_i|A_{1:i-1}) \end{aligned}$$

$i:j$  — sequence notations

- Chain rule applied to sequences of words

Ordered.

$$\begin{aligned} P(w_1, \dots, w_N) &= P(w_1) \cdot P(w_2|w_1) \cdot P(w_3|w_{1:2}) \cdot \dots \cdot P(w_N|w_{1:N-1}) \\ &= \prod_{i=1}^N P(w_i|w_{1:i-1}) \end{aligned}$$

# Probabilities of Sentences

- Calculating the probabilities using Maximum Likelihood Estimations

$$P(w_n | w_{1:n-1}) = \frac{\text{Count}(\overbrace{w_1:n-1}^{\text{prefix}} w_n)}{\sum_w \text{Count}(w_{1:n-1} \underbrace{w}_{\substack{\text{not } w_n \\ \text{Any possible word}}} )} = \frac{\text{Count}(w_{1:n})}{\text{Count}(w_{1:n-1})}$$

*current word*

*word*

*not  $w_n$*

*Any possible word can follow.*

# Probabilities of Sentences — Example

## (1) Application of Chain Rule

$$\begin{aligned} P(\text{"remember to submit your assignment"}) &= P(\text{"remember"}) \cdot \\ &\quad P(\text{"to" | "remember"}) \cdot \\ &\quad P(\text{"submit" | "remember to"}) \cdot \\ &\quad P(\text{"your" | "remember to submit"}) \cdot \\ &\quad P(\text{"assignment" | "remember to submit your"}) \end{aligned}$$

## (2) Maximum Likelihood Estimation

$$P(\text{"remember"}) = \frac{\text{Count}(\text{"remember"})}{N}$$

total  
of words in corpus!

$$P(\text{"to" | "remember"}) = \frac{\text{Count}(\text{"remember to"})}{\text{Count}(\text{"remember"})}$$

...

$$P(\text{"assignment" | "remember to submit your"}) = \frac{\text{Count}(\text{"remember to submit your assignment"})}{\text{Count}(\text{"remember to submit your"})}$$

# Probabilities of Sentences — Problems

$$P(\text{"assignment"} \mid \text{"remember to submit your"}) = \frac{\text{Count}(\text{"remember to submit your assignment"})}{\text{Count}(\text{"remember to submit your"})}$$

- Problem: (very) long sequences

- Large number of entries in table with joint probabilities

- A sequence (or subsequence)  $w_{i:j}$  may not be present in corpus
- $\left. \right\} \rightarrow \text{Count}(w_{i:j}) = 0 \rightarrow \prod_{n=1}^N P(w_n | w_{1:n-1}) = 0$

(we can ignore  $\frac{0}{0}$  here; this can be handled in the implementation)

→ Can we keep the sequences short?

# Natural Language Processing: Foundations

Section 3 — n-Gram Language Models (Markov Assumption)

# Markov Assumption

- Probabilities depend on only on the last  $k$  words

$$P(w_1, \dots, w_N) = \prod_{n=1}^N P(w_n | w_{1:n-1}) \stackrel{\textcolor{violet}{\sim}}{=} \prod_{n=1}^N P(w_n | w_{\textcolor{red}{n-k}:n-1})$$

- For our example:

$$P(\text{"assignment"} | \text{"remember to submit your"}) \approx P(\text{"assignment"} | \text{"your"}) \quad \textcolor{violet}{K=1}$$

$$P(\text{"assignment"} | \text{"submit your"}) \quad \textcolor{violet}{K=2}$$

$$P(\text{"assignment"} | \text{"to submit your"}) \quad \textcolor{violet}{K=3}$$

...

# n-Gram Models

- how many token in gram      context $\rightarrow$
- Unigram (1-gram):**  $P(w_n | w_{1:n-1}) \approx P(w_n)$        $k=0$       context. independent.
- Bigram (2-gram):**  $P(w_n | w_{1:n-1}) \approx P(w_n | w_{n-1})$        $k=1$       } context dependent
- Trigram (3-gram):**  $P(w_n | w_{1:n-1}) \approx P(w_n | w_{n-2}, w_{n-1})$        $k=2$       }
- ⋮

# n-Gram Models

Vocabulary,  $V$

$\Theta(V^1)$

**Unigram** (1-gram):  $P(w_n | w_{1:n-1}) \approx P(w_n)$

$\Theta(V^1)$

**Bigram** (2-gram):  $P(w_n | w_{1:n-1}) \approx P(w_n | w_{n-1})$

$\Theta(V^2)$

**Trigram** (3-gram):  $P(w_n | w_{1:n-1}) \approx P(w_n | w_{n-2}, w_{n-1})$

$\Theta(V^3)$

## Maximum Likelihood Estimation

$$P(w_n) = \frac{\text{Count}(w_n)}{\#\text{words}}$$

$$P(w_n | w_{n-1}) = \frac{\text{Count}(w_{n-1}w_n)}{\text{Count}(w_{n-1})}$$

$$P(w_n | w_{n-1}, w_{n-2}) = \frac{\text{Count}(w_{n-2}w_{n-1}w_n)}{\text{Count}(w_{n-2}w_{n-1})}$$

General MLE for  $n$ -grams:  $P(w_i | w_{n-N+1:n-1}) = \frac{\text{Count}(w_{n-N+1:i})}{\text{Count}(w_{n-N+1:n-1})}$

- **n-Gram models in practice**

- 3-gram, 4-gram, 5-gram models very common
- The larger the n-grams, the more data required

# n-Gram Models — Bigram Example

Example corpus with 3 sentences

<s> I am Sam </s>

<s> Sam I am </s>

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

$$P("I" | " <s>") = \frac{Count(" <s> I")}{Count(" <s>")} =$$

$$P("am" | "I") = \frac{Count("I am")}{Count("I")} =$$

$$P("Sam" | "am") = \frac{Count("am Sam")}{Count("am")} =$$

$$P(" </s>" | "Sam") = \frac{Count("Sam </s>")}{Count("Sam")} =$$

# n-Gram Models — Bigram Example

Example corpus with 3 sentences

<s> I am Sam </s>

<s> Sam I am </s>

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

$$P("I" | " <s>") = \frac{Count(" <s> I")}{Count(" <s> ")} = \frac{2}{3}$$

$$P("am" | "I") = \frac{Count("I am")}{Count("I ")} = \frac{2}{3}$$

$$P("Sam" | "am") = \frac{Count("am Sam")}{Count("am ")} = \frac{1}{2}$$

$$P(" </s>" | "Sam") = \frac{Count("Sam </s>")}{Count("Sam ")} = \frac{1}{2}$$

# n-Gram Models — Bigram Example (25,000 Movie Reviews)

$$P(" < s > \ i \ like \ the \ story \ < /s > ") = ???$$

**Unigram counts:**

i	like	the	story
87,185	19,862	33,0867	11,094

**Bigram counts:**

	i	like	the	story
i	0	693	20	0
like	326	0	1,997	8
the	15	42	0	5,171
story	23	16	16	0

# n-Gram Models — Bigram Example (25,000 Movie Reviews)

$$P(" < s > \ i \ like \ the \ story \ < /s > ") = ???$$

Unigram counts:

i	like	the	story
87,185	19,862	33,0867	11,094

Bigram counts:

	i	like	the	story
i	0	693	20	0
like	326	0	1,997	8
the	15	42	0	5,171
story	23	16	16	0

Bigram probabilities:

	i	like	the	story
i	0.0	0.007949	0.000229	0.0
like	0.016413	0	0.100544	0.000403
the	0.000045	0.000127	0.0	0.015629
story	0.002073	0.001442	0.001442	0.0

Example calculation:

$$P("like" | "i") = \frac{Count("i \ like")}{Count("i")} = \frac{693}{87185} = 0.007949$$

# n-Gram Models — Bigram Example (25,000 Movie Reviews)

**Bigram probabilities:**

	i	like	the	story
i	0.0	<b>0.007949</b>	0.000229	0.0
like	0.016413	0.0	<b>0.100544</b>	0.000403
the	0.000045	0.000127	0.0	<b>0.015629</b>
story	0.002073	0.001442	0.001442	0.0

**Other entries not in the table:**

$$P("i" | " </s>") = 0.088198$$

$$P(" </s>" | "story") = 0.001262$$

$$\begin{aligned} P(" <s> i like the story </s>") &= P("i" | " <s>") \cdot \\ &\quad P("like" | "i") \cdot \\ &\quad P("the" | "like") \cdot \\ &\quad P("story" | "the") \cdot \\ &\quad P(" </s>" | "story") \end{aligned}$$

$$\begin{aligned} P(" <s> i like the story </s>") &= 0.088198 \cdot \\ &\quad 0.007949 \cdot \\ &\quad 0.100544 \cdot \\ &\quad 0.015629 \cdot \\ &\quad 0.001262 \end{aligned}$$

$$P(" <s> i like the story </s>") = 0.00000000139$$

# n-Gram Models — Practical Consideration

- In general

- Each  $P(w_n|w_{1:n-1})$  is rather small →  $\prod_{n=1}^N P(w_n|w_{\textcolor{red}{n-k}:n-1})$  is very small
- Risk of arithmetic underflow

→ Always use an equivalent logarithmic format

- Logarithm is a strictly monotonic function

$$\begin{aligned}P_1 \cdot P_2 \cdot P_3 \cdot \dots \cdot P_N &\propto \log(P_1 \cdot P_2 \cdot P_3 \cdot \dots \cdot P_N) \\&= \log P_1 + \log P_2 + \log P_3 + \dots + \log P_N\end{aligned}$$

Check the options that are correct.

k is for the size of the gram (number of tokens)

n is for the size of the gram (number of tokens)

k is for the size of the context (number of tokens)

n is for the size of the context (number of tokens)



What happens we don't use a special token to mark the end or beginning of a sentence / text?

We will have slightly fewer grams in our vocabulary.

We won't be able to distinguish the start or end of a text or input.

Probability estimates will generally be poorer for short inputs.



# **Natural Language Processing: Foundations**

Section 3 — n-Gram Language Models (Challenges)

# Dealing with Sparsity

- Motivation

- Let's assume we have a bigram language model based 25,000 movie reviews
- Task: Compute the probability of the sentence "*i grasp the story*"  
(compare with the previous example sentence "*i love the story*")

$$\begin{aligned} P(" < s > \ i \ grasp \ the \ story \ < /s > ") = & P("i" | " < s > ") \cdot \\ & P("grasp" | "i") \cdot \\ & P("the" | "grasp") \cdot \\ & P("story" | "the") \cdot \\ & P(" < /s >" | "story") \end{aligned}$$

$$\begin{aligned} P(" < s > \ i \ grasp \ the \ story \ < /s > ") = & 0.088198 \cdot \\ & 0.0 \cdot \leftarrow \\ & 0.226190 \cdot \\ & 0.015629 \cdot \\ & 0.001262 \end{aligned}$$

$$P(" < s > \ i \ grasp \ the \ story \ < /s > ") = 0.0$$

**Problem:** our dataset to train the language model did not contain the bigram "*i grasp*".

$$\text{Count}("i \ grasp") = 0$$

# Dealing with Sparsity

- Problem: **zero-counts**

- n-grams not seen during the training of the language model
- Very common in practice since our training data will be sparse  
(the likelihood for this naturally increases with the size of the n-grams)

$$P(" < s > \ i \ grasp \ the \ story \ </s>") = 0.0$$

**Question:** Is it really acceptable to assign a probability of 0 to a perfectly fine sentence?

- Basic idea: avoid zero-counts

- Assume that any n-gram has a non-zero probability

→ **Smoothing**

# **Natural Language Processing: Foundations**

Section 3 — n-Gram Language Models (Smoothing I: Laplace Smoothing)

# Smoothing

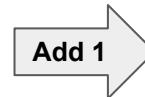
- Basic idea

- Avoid assigning probabilities of 0 to unseen n-grams
- "Move" some probability mass from more frequent n-grams to unseen n-grams
- Also called: **discounting**

- Basic method: **Laplace Smoothing** (also: Add-1 Smoothing)

- Example for bigrams

	i	like	the	story
i	0	693	20	0
like	326	0	1,997	8
the	15	42	0	5,171
story	23	16	16	0



	i	like	the	story
i	1	694	21	1
like	327	1	1,998	9
the	16	43	1	5,172
story	24	17	17	1

no zero counts  
↓

# Smoothing — Laplace Smoothing

- Calculating the probabilities

$$\begin{aligned} P_{\text{Laplace}}(w_n | w_{1:n-1}) &= \frac{\text{Count}_{\text{Laplace}}(w_{1:n-1} w_n)}{\sum_w \text{Count}_{\text{Laplace}}(w_{1:n-1} w)} \\ &= \frac{\text{Count}(w_{1:n-1} w_n) + 1}{\sum_w [\text{Count}(w_{1:n-1} w) + 1]} \\ &= \frac{\text{Count}(w_{1:n-1} w_n) + 1}{\text{Count}(w_{1:n-1}) + V} \end{aligned}$$

vocabulary  
全部詞彙.

e.g., for bigrams:  $P_{\text{Laplace}}(w_n | w_{n-1}) = \frac{\text{Count}(w_{n-1} w_n) + 1}{\text{Count}(w_{n-1}) + V}$

# Smoothing — Laplace Smoothing

- Effects of smoothing on probabilities

Bigram probabilities (without Laplace Smoothing):

	i	like	the	story
i	0.0	0.007949	0.000229	0.0
like	0.016413	0	0.100544	0.000403
the	0.000045	0.000127	0.0	0.015629
story	0.002073	0.001442	0.001442	0.0

Bigram probabilities (with Laplace Smoothing):

	i	like	the	story
i	0.000006	0.004075	0.000123	0.000006
like	0.003175	0.000010	0.019401	0.000087
the	0.000039	0.000104	0.000002	0.012493
story	0.000255	0.000180	0.000180	0.000011

- Observations

- No zero probabilities (duh!)
- Some non-zero probabilities have changed quite a bit!
- For some n-grams: (arguably) too much probability gets moved to zero probabilities

# Smoothing — Laplace Smoothing

- Effects of smoothing on counts

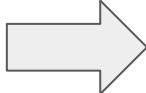
- Question: What counts — without smoothing — would yield  $P_{Laplace}(w_i|w_{i-1})$  ?

$$P_{Laplace}(w_n|w_{n-1}) = \frac{Count(w_{n-1}w_n) + 1}{Count(w_{n-1}) + V} = \frac{Count^*(w_{n-1}w_n)}{Count(w_{n-1})}$$

→  $Count^*(w_{n-1}w_n) = (Count(w_{n-1}w_n) + 1) \cdot \frac{Count(w_{n-1})}{Count(w_{n-1}) + V}$

Bigram counts (original):

	i	like	the	story
i	0	693	20	0
like	326	0	1,997	8
the	15	42	0	5,171
story	23	16	16	0



Bigram counts (adjusted):

	i	like	the	story
i	0.51	355.28	10.75	0.51
like	63.07	0.19	385.34	1.74
the	12.79	34.37	0.80	4133.5
story	2.83	2.00	2.00	0.12

# Add- $k$ Smoothing

- Generalize Laplace (Add-1) Smoothing

- Add  $k$  instead of 1
- Set  $0 < k \leq 1$



$$P_{add-k}(w_n | w_{n-1}) = \frac{Count(w_{n-1}w_n) + k}{Count(w_{n-1}) + kV}$$

$k$  is hyper parameter

# Natural Language Processing: Foundations

Section 3 — n-Gram Language Models (Smoothing II: Common Refinements)

# Backoff & Interpolation

插入、添写。

- Intuition: Utilize less context if required

- Assume we want to calculate  $P(w_n | w_{n-2}, w_{n-1})$  but trigram  $w_{n-2}w_{n-1}w_n$  is not in the dataset

## (1) Backoff

- Make use of bigram probability  $P(w_n | w_{n-1})$  一次不搞那种玩意。
- If still insufficient, use unigram probability  $P(w_n)$

## (2) Interpolation

- Estimate  $P(w_n | w_{n-2}, w_{n-1})$  as a weighted mix of trigram, bigram, and unigram probabilities
- Learn weights  $\lambda_i$  from data
- In practice better than Backoff

$$\text{Count}(w_{n-2} w_{n-1} w_n) = 0$$



# Linear Interpolation (example for trigrams)

- Simple interpolation

$$\hat{P}(w_n | w_{n-2}, w_{n-1}) = \lambda_1 P(w_n) + \lambda_2 P(w_n | w_{n-1}) + \lambda_3 P(w_n | w_{n-2}, w_{n-1})$$

*independent from the extra bigram*

*independent unigram*

with  $\sum_i \lambda_i = 1$

- $\lambda_i$  conditional on context

$$\hat{P}(w_n | w_{n-2}, w_{n-1}) = \lambda_1(w_{n-2}, w_{n-1}) P(w_n) + \lambda_2(w_{n-2}, w_{n-1}) P(w_n | w_{n-1}) + \lambda_3(w_{n-2}, w_{n-1}) P(w_n | w_{n-2}, w_{n-1})$$

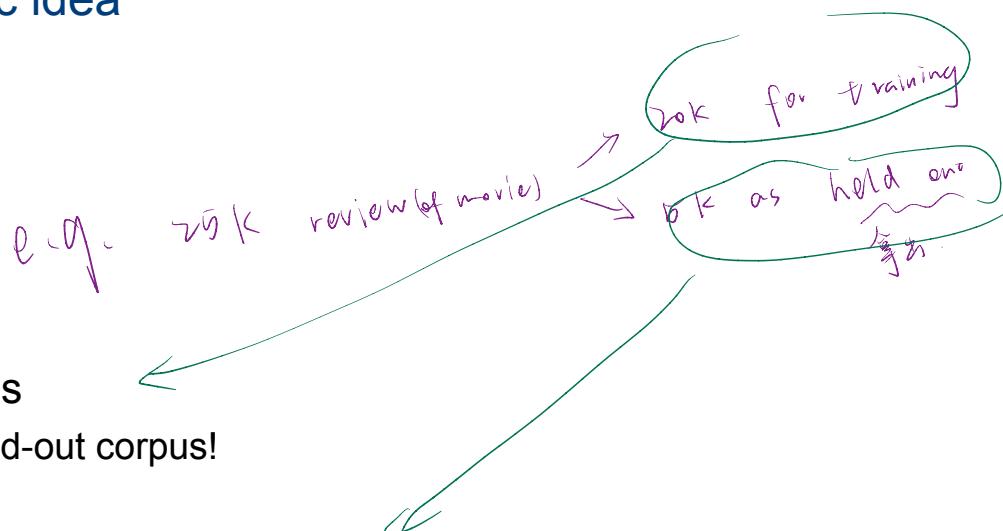
*dependent or the extra bigram*

# Backoff & Interpolation

- Learn weights  $\lambda_i$  from data — basic idea

- (1) Collect held-out corpus

- Additional corpus *or*
- Split from initial corpus



- (2) Calculate all n-gram probabilities

- Calculation must no consider held-out corpus!

- (3) Find  $\lambda_i$  that maximize  $\hat{P}(w_n | w_{n-2}, w_{n-1})$  over held-out corpus

- e.g., using Expectation-Maximization (EM) algorithm (not further discusses here)

# **Natural Language Processing: Foundations**

Section 3 — n-Gram Language Models (Smoothing III: Advanced Techniques)

# Kneser-Ney Smoothing

most common model

- Idea of Kneser-Ney Smoothing: **Absolute Discounting Interpolation**

Remove a fixed value  
from all bigram counts

Interpolation but with better  
estimates for unigram probabilities

$$P_{KN}(w_n|w_{n-1}) = \frac{\max [Count(w_{n-1}w_n) - d, 0]}{Count(w_{n-1})} + \lambda(w_{n-1}) P_{KN}(w_n)$$

0.75  
詳情看 Stanford · PPT.

**Note:** We only look at a bigram language model in the following to keep the examples and notations easy. Kneser-Ney Smoothing analogously defined for larger n-grams.

# Kneser-Ney Smoothing — Absolute Discounting

- **Absolute discounting**

- Remove fixed value  $d$  from bigram counts  
(typically:  $0 < d < 1$ )
- Makes probability mass for unigrams available
- Intuition

If  $\text{Count}(w_{n-1}w_n)$  is large, count hardly affected

If  $\text{Count}(w_{n-1}w_n)$  is small, count not that useful to begin with

just a fail-safe to avoid negative probabilities

$$\frac{\max[\text{Count}(w_{n-1}w_n) - d, 0]}{\text{Count}(w_{n-1})}$$

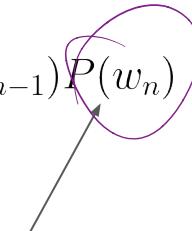
→ Question: How to pick the value(s) for  $d$  ?

0.75

# Kneser-Ney Smoothing — Interpolation with a Twist

- Motivation

$$P_{KN}(w_n|w_{n-1}) = \frac{\max [Count(w_{n-1}w_n) - d, 0]}{Count(w_{n-1})} + \lambda(w_{n-1}) P(w_n)$$



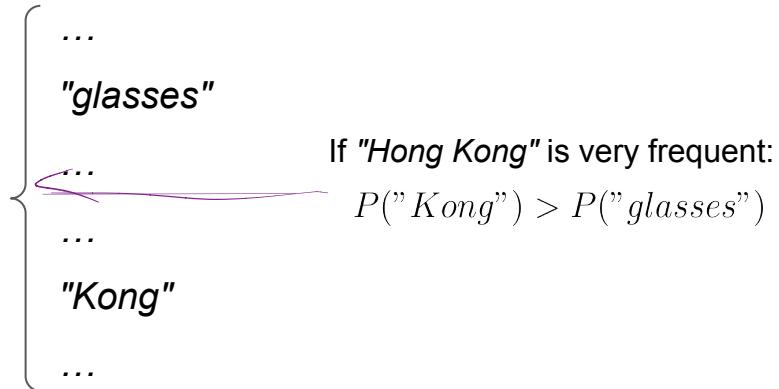
Using basic interpolation, that would just be the unigram probability

→ But is this actually a good idea?

Predict the missing word:

"I can't see without my reading Kong"

↑  
bad choice.



# Kneser-Ney Smoothing — Interpolation with a Twist

- The difference between "glasses" and "Kong" — Intuition

- "glasses" is preceded by many other words
- "Kong" is almost only preceded by "Hong"

↑  
香港.

→  $P(w) = \text{"How likely is } w \text{ ?"}$  maybe not most intuitive approach

- Alternative:  $P_{KN}(w) = \text{"How likely is } w \text{ to appear as a novel continuation?"}$

- $P_{KN}(w)$  is high  $\Leftrightarrow$  there are many words  $w'$  that form an existing bigram  $w'w \Rightarrow \text{glass}$
- $P_{KN}(w)$  is low  $\Leftrightarrow$  there are only few words  $w'$  that form an existing bigram  $w'w \Rightarrow (\text{Hong})\text{Kong}.$

→ How can we quantify this?

# Kneser-Ney Smoothing — Interpolation with a Twist

- Calculating  $P_{KN}(w)$

↳ using counts,

reading  
w ↗ sun  
glasses.  
glasses.

#words  $w'$  that form an existing bigram  $w'w$

$$P_{KN}(w) = \frac{|\{w' : \text{Count}(w'w) > 0\}|}{|\{(u, v) : \text{Count}(uv) > 0\}|}$$

total number of existing bigrams

normalization to ensure that

$$\sum_{n=1}^N P(w_n) = 1$$

# Kneser-Ney Smoothing — Wrapping it Up

$$P_{KN}(w_n|w_{n-1}) = \frac{\max [Count(w_{n-1}w_n) - d, 0]}{Count(w_{n-1})} + \lambda(w_{n-1}) P_{KN}(w_n)$$

↓  
 $\sum P_{KN} = 1$       ↓

*new unigram probabilities,*  
last missing puzzle piece

- Normalizing factor  $\lambda$

- Required to account for the probability mass we have discounted

$$\lambda(w_{n-1}) = \underbrace{\frac{d}{Count(w_{n-1})}}_{\text{normalized discount}} \cdot \underbrace{|\{w' : Count(w_{n-1}w') > 0\}|}_{\#\text{words that can follow}}$$

= #words that have been discounted  
= #times the normalized discount has been applied

# **Natural Language Processing: Foundations**

Section 3 — n-Gram Language Models (Evaluation)

# Evaluating Language Models

用 范文 比较 的 文章 对练  
不 符 合 现代 人 .

- A Language Model (LM) is considered good if
  - It assigns high probabilities to frequently occurring sentences in a corpus
  - It assigns low probabilities otherwise (e.g., rarely occurring sentences)
- 2 basic approaches to compare LMs

## 外部 Extrinsic Evaluation

- Requires a downstream task (e.g., spell checker, speech recognition)
- Run a downstream task with each LM and compare the results
- Can be very expensive and time-consuming

## Intrinsic Evaluation

- Evaluate each LM on a test corpus
  - Generally cheaper and faster
  - Require an intrinsic metric to compare LMs
- **Perplexity** (among other metrics)

# Intrinsic Evaluation

- 3 core steps for an intrinsic evaluation
  - (1) Train the LM on a **training corpus**  
(i.e., compute the n-gram probabilities)
  - (2) Tune parameters of the LM using a **development corpus**  
(e.g.,  $k$  in case of Add- $k$  Smoothing,  $d$  in case of Kneser-Ney Smoothing)
  - (3) Compute evaluation metrics on a **test corpus**  
(e.g., perplexity)
- Common corpus breakdown: 80/10/10 (80% training, 10% development, 10% test)

# Perplexity

- **Perplexity — Definition**

- Inverse probability of test corpus  $W$
- Normalized by the number of words  $N$  in test corpus

$$PP(W) = P(w_1, w_2, \dots, w_N)^{-\frac{1}{N}}$$

geometric mean.

$$= \sqrt[N]{\frac{1}{P(w_1, w_2, \dots, w_N)}}$$

chain rule:

$$= \sqrt[N]{\prod_{n=1}^N \frac{1}{P(w_n|w_1, \dots, w_{n-1})}}$$

逐字取次.

e.g., for bigrams:

$$= \sqrt[N]{\prod_{n=1}^N \frac{1}{P(w_n|w_{n-1})}}$$

Minimizing perplexity  $\Leftrightarrow$  Maximizing probability

Perplexity, 翻译为中文的意思为：困惑、混乱。在NLP中，表示 语言模型 的困惑度，是交叉熵的指数形式，可以作为语言模型的评价指标。Perplexity越低，说明模型拟合效果越好。Perplexity 计算公式如下：

$$\begin{aligned} \text{Perplexity}(S) &= p(w_1, w_2, w_3, \dots, w_m)^{-1/m} \\ &= \sqrt[m]{\prod_{i=2}^m \frac{1}{p(w_i | w_1, w_2, w_3, \dots, w_{i-1})}} \end{aligned}$$

简单来说，perplexity刻画的是语言模型预测一个语言样本的能力，比如已经知道了( $w_1, w_2, w_3, \dots, w_m$ )这句话会出现在语料库之中，那么通过语言模型计算得到这句话的概率越高，说明语言模型对这个语料库拟合的越好。

在语言模型的训练中，通常采用perplexity的对数表达形式：

$$\log(\text{perplexity}(S)) = -\frac{1}{m} \sum_{i=2}^m p(w_i, w_2, w_3, \dots, w_m)$$

在数学上， $\log$  perplexity 可以看作真实分布与预测分布之间的 交叉熵 Cross Entropy，交叉熵描述了两个概率分布之间的一种距离，假设  $x$  是一个离散变量， $u(x), v(x)$  是两个与  $x$  相关的概率分布，那么  $u, v$  之间的交叉熵的定义是分布  $u$  下  $-\log(v(x))$  的期望值：

$$H(u, v) = E_u[-\log(v(x))] = - \sum_x u(x) \log(v(x))$$

# Perplexity — Intuition

- When is the perplexity **high**?

Many n-grams are frequent in the training corpus but rare in the test corpus



Very few high  $P(w_n|w_{n-1})$  values over test corpus



High perplexity  $PP(W) = \sqrt[N]{\prod_{n=1}^N \frac{1}{P(w_n|w_{n-1})}}$

Many n-grams are rare in the training corpus but frequent in the test corpus



Many low  $P(w_n|w_{n-1})$  values over test corpus



# Perplexity — Practical Consideration

- In general
  - Each  $P(w_n|w_{1:n-1})$  rather small  $\rightarrow \prod_{n=1}^N P(w_n|w_{1:n-1})$  very small
  - Risk of arithmetic underflow

- Again, logarithm to the rescue

$$\ln PP(W) = -\frac{1}{N} P(w_1, w_2, \dots, w_N)$$

$$PP(W) = e^{\ln PP(W)}$$

$$= -\frac{1}{N} \ln \prod_{n=1}^N P(w_n|w_1, \dots, w_{n-1})$$

$$= -\frac{1}{N} \sum_{n=1}^N \ln P(w_n|w_1, \dots, w_{n-1})$$

e.g., for bigrams:

$$= -\frac{1}{N} \sum_{n=1}^N \ln P(w_n|w_{n-1})$$

# Perplexity — Toy Example

- Evaluation setup

- Bigram LM trained over 25k movie reviews
- Small test corpus  $W$  with  $N = 12$

$$W = \begin{bmatrix} "⟨s⟩ i like good movies ⟨/s⟩", \\ "⟨s⟩ the story is funny ⟨/s⟩" \end{bmatrix}$$

$$PP(W) = \sqrt[N]{\prod_{n=1}^N \frac{1}{P(w_n|w_{n-1})}} = 40.1$$

bigram	P(bigram)
"⟨s⟩ i"	0.0882
"i like"	0.0079
"like good"	0.0013
"good movies"	0.0062
"movies ⟨/s⟩"	0.0034
"⟨s⟩ the"	0.0990
"the story"	0.0156
"story is"	0.1138
"is funny"	0.0022
"funny ⟨/s⟩"	0.0081

# Perplexity — Real-World Example

- Evaluation setup

- Unigram, Bigram, Trigram LMs trained over *Wall Street Journal* articles
- Training corpus: ~38 million words (~20k unique words)
- Test corpus: ~1.5 million words

	Unigram	Bigram	Trigram
Perplexity	962	170	109

越低 越好 .

# Summary

- Language Models — assigning probabilities to sentences
  - Very important concept for many NLP tasks
  - Different methods to compute sentence probabilities  
(here: n-grams; later we come back to them using neural networks)
- n-gram Language Models
  - Intuitive training → Maximum Likelihood Estimations
  - Main consideration: zero probabilities due to large n-grams and/or open vocabularies

Markov Assumption to limited size of considered n-grams

Focus here: Smoothing  
(maybe with backoff & interpolation)



In practice, typically a combination of these and similar approaches

Step 1:

add padding :

把 [ "]", "am", "Bob" ], [ ... ], ... ]

变成 [ "<s>", "]", "am", "Bob", "</s>" ], ... ]

padding\_text = [ ] , sos = '<s>' , eos = '</s>'

for text in tokenized\_texts:

padded\_text = [ self.sos ] + text + [ self.eos ]

padding\_texts.append(padded\_text),

step 2. build vocabulary

vocabulary = set()