

Language N-Gram Model

Language N-Gram Model– Module II

Probabilistic Language Models

- Today's goal: assign a probability to a sentence

- Machine Translation:

- $P(\text{high winds tonite}) > P(\text{large winds tonite})$

- Spell Correction

- The office is about fifteen **minuets** from my house

- $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$

- Speech Recognition

- $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$

- + Summarization, question-answering, etc., etc.!!

Why?

Probabilistic Language Modeling

- Goal: compute the probability of a sentence or sequence of words:

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

- Related task: probability of an upcoming word:

$$P(w_5 | w_1, w_2, w_3, w_4)$$

- A model that computes either of these:

$P(W)$ or $P(w_n | w_1, w_2 \dots w_{n-1})$ is called a **language model**.

- Better: **the grammar** But **language model** or **LM** is standard

What is a language model?

- Probability distributions over sentences (i.e., word sequences) : $\mathbf{P(W) = P(w_1 w_2 w_3 w_4 \dots w_k)}$

- Can use them to generate strings

- $\mathbf{P(w_k \mid w_2 w_3 w_4 \dots w_{k-1})}$

- Rank possible sentences

- $\mathbf{P(\text{“Today is Thursday”}) > P(\text{“Thursday Today is”})}$

- $\mathbf{P(\text{“Today is Thursday”}) > P(\text{“Today is Sunny”})}$

Uses of Language Models

- Speech recognition

- “I ate a cherry” is a more likely sentence than “Eye eight uh Jerry”

- OCR & Handwriting recognition

- More probable sentences are more likely correct readings.

- Machine translation

- More likely sentences are probably better translations.

- Generation

- More likely sentences are probably better NL generations.

- Context sensitive spelling correction

- “There are problems with this sentence.”

Language model applications

- Context-sensitive spelling correction

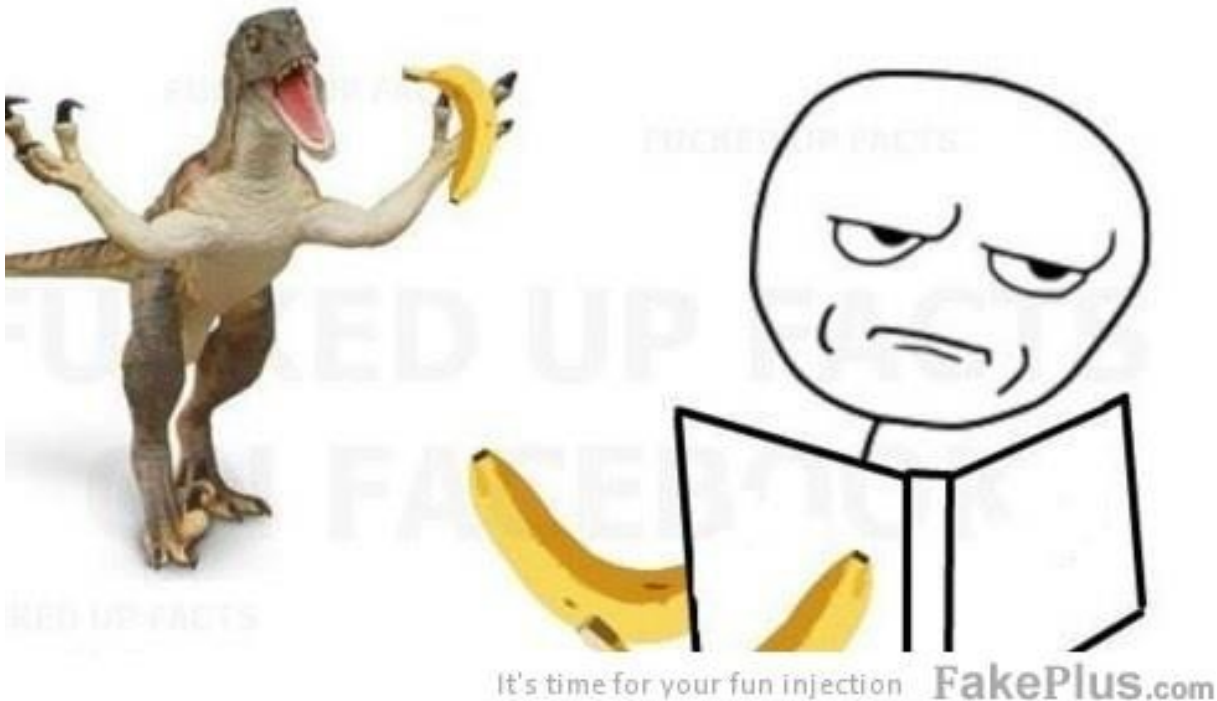
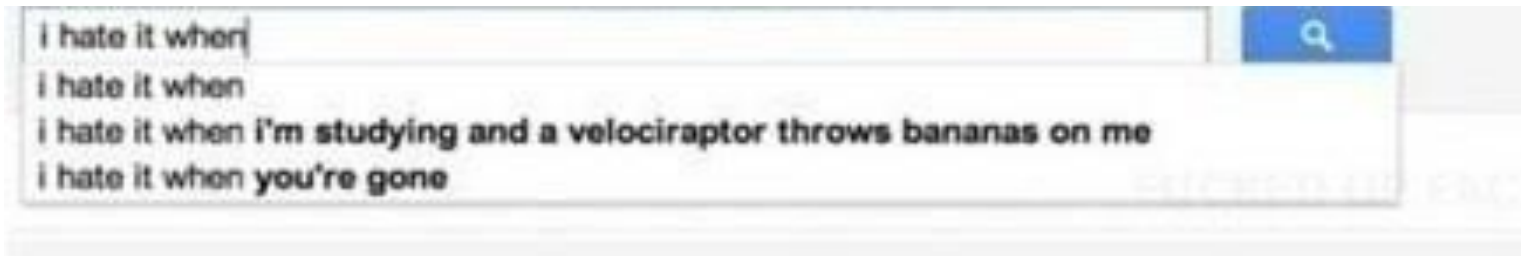


Real Word Spelling Errors

- Words that sound the same
 - Their/they're/there
 - To/too/two
 - Weather/whether
 - Peace/piece
 - You're/your
- Typos that result in real words
 - Lave for Have

Language model applications

Autocomplete



Language model applications

- **Language generation :**
<https://pdos.csail.mit.edu/archive/scigen/>

Deploying Superblocks and Compilers

Julia and Dan

Abstract

Recent advances in replicated algorithms and relational symmetries have paved the way for architecture. After years of natural research into erasure coding, we show the deployment of courseware, which embodies the key principles of steganography. *Loy*, our new system for the exploration of sensor networks, is the solution to all of these issues.

1 Introduction

Steganographers agree that robust symmetries are an interesting new topic in the field of cryptography, and information theorists concur. We view operat-

thesize unstable algorithms, we fulfil without investigating the evaluation of

Our contributions are threefold. First, how erasure coding can be applied to reinforcement learning. We present an algorithm for the deployment of extremizing (*Loy*), which we use to prove that and operating systems [19, 7, 14] can fill this goal. We examine how replication can be applied to the deployment of linked

The rest of this paper is organized as follows. First, we motivate the need for fiber. We demonstrate the synthesis of the T. Finally, we conclude.

Completion Prediction

- A language model also supports predicting the completion of a sentence.
 - Please turn off your cell
 - Kindly submit
- *Predictive text input* systems can guess what you are typing and give choices on how to complete it.

N-Gram Models of Language

- Use the previous $N-1$ words in a sequence to predict the next word
- Language Model (LM)
 - unigrams, bigrams, trigrams,...
- How do we train these models?
 - Very large corpora

Corpora

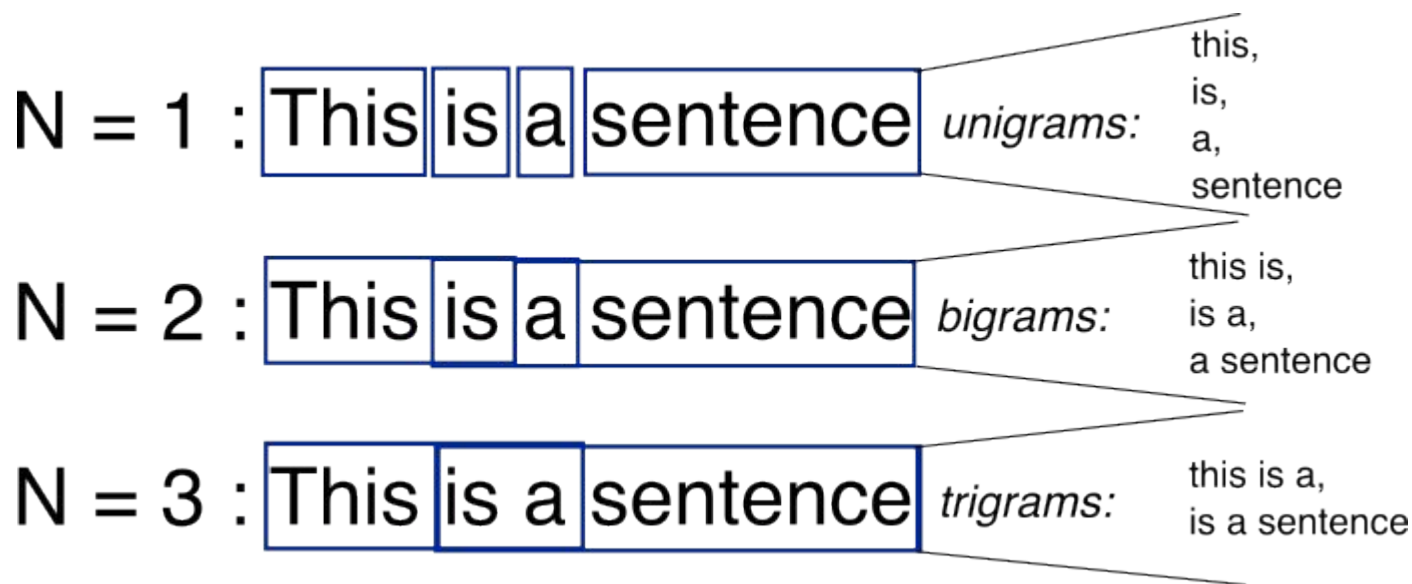
- **Corpora are online collections of text and speech**
 - Brown Corpus
 - Wall Street Journal
 - AP newswire
 - Hansards
 - Timit
 - DARPA/NIST text/speech corpora (Call Home, Call Friend, ATIS, Switchboard, Broadcast News, Broadcast Conversation, TDT, Communicator)
 - TRAINS, Boston Radio News Corpus

Terminology

- **Sentence:** unit of written language
- **Utterance:** unit of spoken language
- **Word Form:** the inflected form as it actually appears in the corpus
- **Lemma:** an abstract form, shared by word forms having the same **stem**, part of speech, word sense – stands for the class of words with same **stem**
- **Types:** number of distinct words in a corpus (vocabulary size)
- **Tokens:** total number of words

Bag-of-Words with N-grams

- **N-grams**: a contiguous sequence of n tokens from a given piece of text



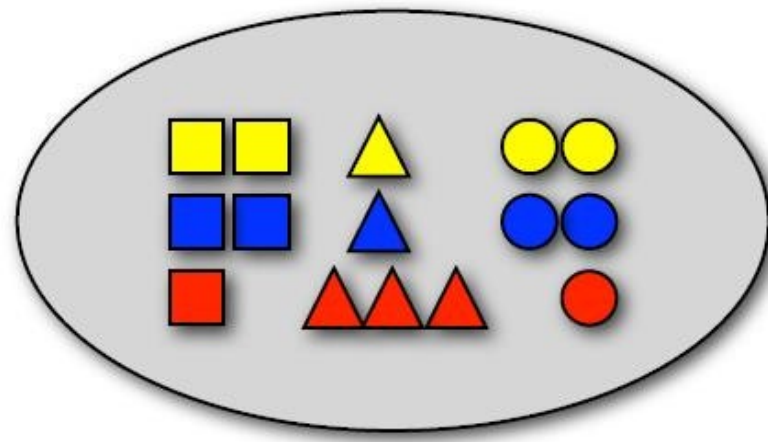
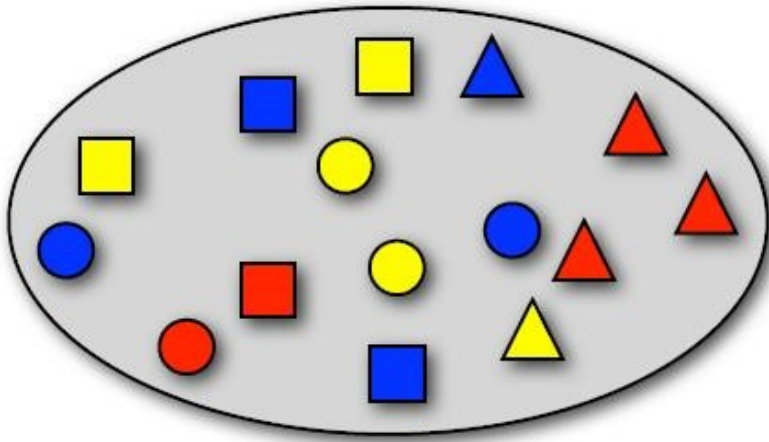
N-Gram Models

- **Unigram model:** $P(w_1)P(w_2) \dots P(w_n)$
- **Bigram model:** $P(w_1)P(w_2|w_1)P(w_3|w_2) \dots P(w_n|w_{n-1})$
- **Trigram model:**
 - $P(w_1)P(w_2|w_1)P(w_3|w_2, w_1) \dots P(w_n|w_{n-1}w_{n-2})$
- **N-gram model:**
 - $P(w_1)P(w_2|w_1) \dots P(w_n|w_{n-1}w_{n-2} \dots w_{n-N})$

Random language via n-gram

- <http://www.cs.jhu.edu/~jason/465/PowerPoint/lect01,3tr-ngram-gen.pdf>
- Behind the scenes – probability theory

Sampling with replacement



1. $P(\blacksquare) = ?$ 2. $P(\square) = ?$ 3. $P(\text{red}, \square) = ?$

4. $P(\text{blue}) = ?$ 5. $P(\text{red} \mid \square) = ?$

6. $P(\square \mid \text{red}) = ?$ 7. $P(\text{red circle, yellow triangle, blue triangle, blue square}) = ?$

8. $P(\square\square\square) = ?$ 9. $P(2 \times \text{red square}, 1 \times \text{blue square}, 4 \times \text{yellow triangle}) = ?$

Sampling words with replacement

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

$$P(\text{of}) = 3/66$$

$$P(\text{Alice}) = 2/66$$

$$P(\text{was}) = 2/66$$

$$P(\text{to}) = 2/66$$

$$P(\text{her}) = 2/66$$

$$P(\text{sister}) = 2/66$$

$$P(,) = 4/66$$

$$P(') = 4/66$$

How to compute $P(W)$?

- How to compute this joint probability:
 - $P(\text{its, water, is, so, transparent, that})$
- Intuition: let's rely on the Chain Rule of Probability

Reminder : Chain Rule

- Recall the definition of conditional probabilities

Rewriting:

- More variables:

$$P(A,B,C,D) = P(A)P(B|A)P(C|A,B)P(D|A,B,C)$$

- The Chain Rule in General

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2) \dots P(x_n|x_1, \dots, x_{n-1})$$

Reminder : Chain Rule

- The big red dog
- $P(\text{The}) * P(\text{big}|\text{the}) * P(\text{red}|\text{the big}) * P(\text{dog}|\text{the big red})$
- Better $P(\text{The} | \text{<Beginning of sentence>})$ written as
 $P(\text{The} | \text{<S>})$

Reminder : Chain Rule

- The **Chain Rule** applied to compute joint probability of words in sentence.

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i \mid w_1 w_2 \dots w_{i-1})$$

$P(\text{"its water is so transparent"}) =$

$P(\text{its}) \times P(\text{water} \mid \text{its}) \times P(\text{is} \mid \text{its water})$

$\times P(\text{so} \mid \text{its water is}) \times P(\text{transparent} \mid \text{its water is so})$

How to estimate these Probabilities ?

- Could we just count and divide?

$$P(\text{the } l \text{ its water is so transparent that}) = \frac{\textit{Count}(\text{its water is so transparent that the})}{\textit{Count}(\text{its water is so transparent that})}$$

- No! Too many possible sentences!
- We'll never see enough data for estimating these

Markov Assumption



Andrei Markov

- Simplifying assumption:

$$P(\text{the l its water is so transparent that}) \approx P(\text{the l that})$$

- Or maybe

$$P(\text{the l its water is so transparent that}) \approx P(\text{the l transparent that})$$

Language model with N-gram

- The chain rule:
 - $$P(X_1, X_2, \dots, X_n) = P(X_1) P(X_2|X_1) P(X_3|X_2, X_1) \dots P(X_n|X_1, \dots, X_{n-1})$$
- N-gram language model assumes each word depends only on the last n-1 words (Markov assumption)

Language model with N-gram

- Example: trigram (3-gram)

$$P(w_n | w_1, \dots, w_{n-1}) \neq P(w_n | w_{n-2}, w_{n-1})$$

$$P(w_1, \dots, w_n) =$$

$$P(w_1)P(w_2 | w_1) \dots P(w_n | w_{n-2}, w_{n-1})$$

- $P(\text{"Today is a sunny day"})$

$$= P(\text{"Today"})P(\text{"is"} | \text{"Today"})P(\text{"a"} | \text{"is"}, \text{"Today"}) \dots$$

$$P(\text{"day"} | \text{"sunny"}, \text{"a"})$$

N-grams : Example - The big red dog

- Unigrams: $P(\text{dog})$
- Bigrams: $P(\text{dog}|\text{red})$
- Trigrams: $P(\text{dog}|\text{big red})$
- Four-grams: $P(\text{dog}|\text{the big red})$

In general, we'll be dealing with
 $P(\text{Word} | \text{Some fixed prefix})$

Language Modeling

Estimating N-gram
Probabilities

Maximum Likelihood Estimation or MLE.

- An intuitive way to estimate probabilities is called **maximum likelihood estimation or MLE**.
- We get maximum likelihood estimation the MLE estimate for the parameters of an n-gram model by getting counts from a corpus, and normalizing the counts so that they lie between 0 and 1
- For example, to compute a particular bigram probability of a word w_n given a previous word w_{n-1} , we'll compute the count of the bigram $C(w_{n-1}w_n)$ and normalize by the sum of all the bigrams that share the same first word w_{n-1} :

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{\sum_w C(w_{n-1}w)}$$

Estimating bigram

The Maximum Likelihood
Estimate

$$P(w_i | w_{1:i-1}) = \frac{\text{count}(w_{1:i-1}, w_i)}{\text{count}(w_{1:i-1})}$$

$$P(w_i | w_{1:i-1}) = \frac{c(w_{1:i-1}, w_i)}{c(w_{1:i-1})}$$

<s> I am Sam </s>

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

<s> Sam I am </s>

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

$$P(I | <s>) = \frac{2}{3} = .67 \quad P(\text{Sam} | <s>) = \frac{1}{3} = .33 \quad P(\text{am} | I) = \frac{2}{3} = .67$$

$$P(</s> | \text{Sam}) = \frac{1}{2} = 0.5 \quad P(\text{Sam} | \text{am}) = \frac{1}{2} = .5 \quad P(\text{do} | I) = \frac{1}{3} = .33$$

We estimate the n-gram probability by dividing the observed frequency of a particular sequence by the observed frequency of a prefix. This ratio is called a relative frequency

More examples: Berkeley Restaurant Project

sentences

can you tell me about any good cantonese restaurants close
by
mid priced thai food is what i'm looking

for tell me about chez panisse

can you give me a listing of the kinds of food that are available
i'm looking for a good place to eat

breakfast when is caffe venezia open

during the day

Bigram counts for eight of the words (out of V corpus of 9332 sentences. Zero counts are in blue

= (1 - 446) in the Berkeley Restaurant Project

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Raw bigram

Normalize by unigrams:

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

Result:

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Bigram estimates of sentence probabilities

$P(<s> \text{ I want english food } </s>) =$

$P(\text{I} | <s>)$

$\times P(\text{want} | \text{I})$

$\times P(\text{english} | \text{want})$

$\times P(\text{food} | \text{english})$

$\times P(</s> | \text{food})$

$= .000031$

What kinds of

$$P(\text{english} | \text{want}) = .0011$$

$$P(\text{chinese} | \text{want}) = .0065$$

$$P(\text{to} | \text{want}) = .66$$

$$P(\text{eat} | \text{to}) = .28$$

$$P(\text{food} | \text{to}) = 0$$

$$P(\text{want} | \text{spend}) =$$

$$0 \quad P(i | \langle s \rangle) = .25$$

bigrams in our tongue twister

Peter Piper picked a peck of pickled pepper.

Where's the pickled pepper that Peter Piper picked?

the conditional probability of “Piper” given
“Peter”:

$$) \ p(\text{Piper}|\text{Peter}) = \frac{|\text{Peter Piper}|}{|\text{Peter}|} = \frac{2}{2} = 1$$

$$p(\text{Piper}|\text{a}) = \frac{|\text{a Piper}|}{|\text{a}|} = \frac{0}{1} = 0$$

Bigrams	Bigram frequencies
picked a	1
pepper that	1
peck of	1
a peck	1
pickled pepper	2
Where s	1
Piper picked	2
the pickled	1
Peter Piper	2
of pickled	1
pepper Where	1
that Peter	1
s the	1

<s> Peter Piper picked a peck of pickled pepper. </s>

<s> Where's the pickled pepper that Peter Piper picked? </s>

$$\begin{aligned} & p\left(\begin{array}{l} \text{Where's the pickled pepper that Peter} \\ \text{Piper picked a peck of pickled pepper} \end{array}\right) \\ = & p(\text{Where's}|\text{<s>}) \times p(\text{the}|\text{Where's}) \times p(\text{pickled}|\text{the}) \times \\ & p(\text{pepper}|\text{pickled}) \times p(\text{that}|\text{pepper}) \times p(\text{Peter}|\text{that}) \times \\ & p(\text{Piper}|\text{Peter}) \times p(\text{picked}|\text{Piper}) \times p(\text{a}|\text{picked}) \times \\ & p(\text{peck}|\text{a}) \times p(\text{of}|\text{peck}) \times p(\text{pickled}|\text{of}) \times \\ & p(\text{peper}|\text{pickled}) \times p(\text{</s>}|\text{pepper}) \\ = & .5 \times 1 \times 1 \times 1 \times .5 \times 1 \times 1 \times 1 \times .5 \times 1 \times 1 \times 1 \times 1 \times .5 \\ = & .0625 \end{aligned}$$

Practical

We do everything in log space

- Avoid underflow
- (also adding is faster than multiplying)

$$\log(p_1 \times p_2 \times p_3 \times p_4) = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

Language Modeling

Tools

SRILM

- <http://www.speech.sri.com/projects/srilm/>

KenLM

- <https://kheafield.com/code/kenlm/>

Google Book

N

<http://ngrams.googlelabs.com/>

Example from Daniel Martin

I want chinese food.

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Figure 3.1 Bigram counts for eight of the words (out of $V = 1446$) in the Berkeley Restaurant Project corpus of 9332 sentences. Zero counts are in gray.

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Figure 3.2 Bigram probabilities for eight words in the Berkeley Restaurant Project corpus of 9332 sentences. Zero probabilities are in gray.

unigram probabilities):

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

Here are a few other useful probabilities:

$$P(i|<s>) = 0.25 \quad P(\text{english}|\text{want}) = 0.0011$$

$$P(\text{food}|\text{english}) = 0.5 \quad P(</s>|\text{food}) = 0.68$$

$$\begin{aligned}
 & i \text{ want english food } </s> \\
 & = P(i|<s>)P(\text{want}|i)P(\text{english}|\text{want}) \\
 & \quad P(\text{food}|\text{english})P(</s>|\text{food}) \\
 & = .25 \times .33 \times .0011 \times 0.5 \times 0.68 \\
 & = .000031
 \end{aligned}$$

0= In the given data set these two words are not coming back to back it is not always true . We need solution

Laplace smoothing(add one)

To keep a language model from assigning zero probability to these unseen events, we'll have to shave off a bit of probability mass from some more frequent events and give it to the events we've never seen.

This modification is called **smoothing or discounting**

For add-one smoothed bigram counts, we need to augment the unigram count by the number of total word types in the vocabulary V :

$$P_{\text{Laplace}}^*(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{\sum_w (C(w_{n-1}w) + 1)} = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + V} \quad (3.23)$$

Thus, each of the unigram counts given in the previous section will need to be augmented by $V = 1446$. The result is the smoothed bigram probabilities in Fig. 3.6.

	i	want	to	eat	chinese	food	lunch	spend
i	0.0015	0.21	0.00025	0.0025	0.00025	0.00025	0.00025	0.00075
want	0.0013	0.00042	0.26	0.00084	0.0029	0.0029	0.0025	0.00084
to	0.00078	0.00026	0.0013	0.18	0.00078	0.00026	0.0018	0.055
eat	0.00046	0.00046	0.0014	0.00046	0.0078	0.0014	0.02	0.00046
chinese	0.0012	0.00062	0.00062	0.00062	0.00062	0.052	0.0012	0.00062
food	0.0063	0.00039	0.0063	0.00039	0.00079	0.002	0.00039	0.00039
lunch	0.0017	0.00056	0.00056	0.00056	0.00056	0.0011	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

Figure 3.6 Add-one smoothed bigram probabilities for eight of the words (out of $V = 1446$) in the BeRP corpus of 9332 sentences. Previously-zero probabilities are in gray.

Problem with add-one smoothing

- Data from the AP from (Church and Gale, 1991)
 - Corpus of 44,000,000 bigram tokens, 22,000,000 for training
 - Vocabulary of 273,266 words, i.e. 74,674,306,760 possible bigrams
 - 74,671,100,000 bigrams were unseen
 - frequency is the number of occurrences per 22,000,000 samples
 - To get probability, divide frequency by 22,000,000
 - each unseen bigram was given a frequency of 0.000295

num. of times
appeared in
training corpus

Freq. observed
in testing
corpus

f_{MLE}	$f_{empirical}$	$f_{add-one}$
0	0.000027	0.000295
1	0.448	0.000589
2	1.25	0.000884
3	2.24	0.001180
4	3.23	0.001470
5	4.21	0.001770

Add-one smoothed
freq. given to
testing corpus

too high

too low

Adding a little of probability over a huge number of unseen events gives too much probability mass to all unseen events
Instead of giving small portion of probability to unseen events, most of the probability space is given to unseen events

Add K

If $k=0.5$, Lidstone's law is called Expected Likelihood estimation or Jeffrey's Perks law.

- instead of adding 1, add some other (smaller) positive value δ

$$P_{\text{Add}\delta}(w_1 w_2 \dots w_n) = \frac{C(w_1 w_2 \dots w_n) + \delta}{N + \delta B}$$

- most widely used value for $\delta = 0.5$
- if $\delta = 0.5$, Lidstone's Law is called:
 - the **Expected Likelihood Estimation** (ELE)
 - or the **Jeffreys-Perks** Law

$$P_{\text{ELE}}(w_1 w_2 \dots w_n) = \frac{C(w_1 w_2 \dots w_n) + 0.5}{N + 0.5 B}$$

- better than add-one, but still not very good

Smoothing: Good Turing



- Imagine you are fishing
 - You have bass, carp, cod, tuna, trout, salmon, eel, shark, tilapia, etc. in the sea
- You have caught 10 Carp, 3 Cod, 2 tuna, 1 trout, 1 salmon, 1 eel
- How likely is it that next species is new?
 - roughly $3/18$, since 18 fish total, 3 unique species
- How likely is it that next is tuna? Less than $2/18$
 - 2 out of 18 are tuna, but we have to give some “room” to the new species that we may catch in the future
- Say that there are 20 species of fish that we have not seen yet (bass, shark, tilapia,....)
- The probability of any individual unseen species is $\frac{3}{18 \cdot 20}$
- $P(\text{shark}) = P(\text{tilapia}) = \frac{3}{18 \cdot 20}$

Smoothing: Good Turing



- How many species (n-grams) were seen once?
 - Let N_1 be the number species (n-grams) seen once
- Use it to estimate for probability of unseen species
 - Probability of new species (new n-gram) is N_1/N
- Let N_0 be the number of unseen species (unseen n-grams). Spreading around the mass equally for unseen n-grams, the probability of seeing any individual unseen species (unseen n-gram) is

$$\frac{N_1}{N \cdot N_0}$$

Smoothing: Good Turing



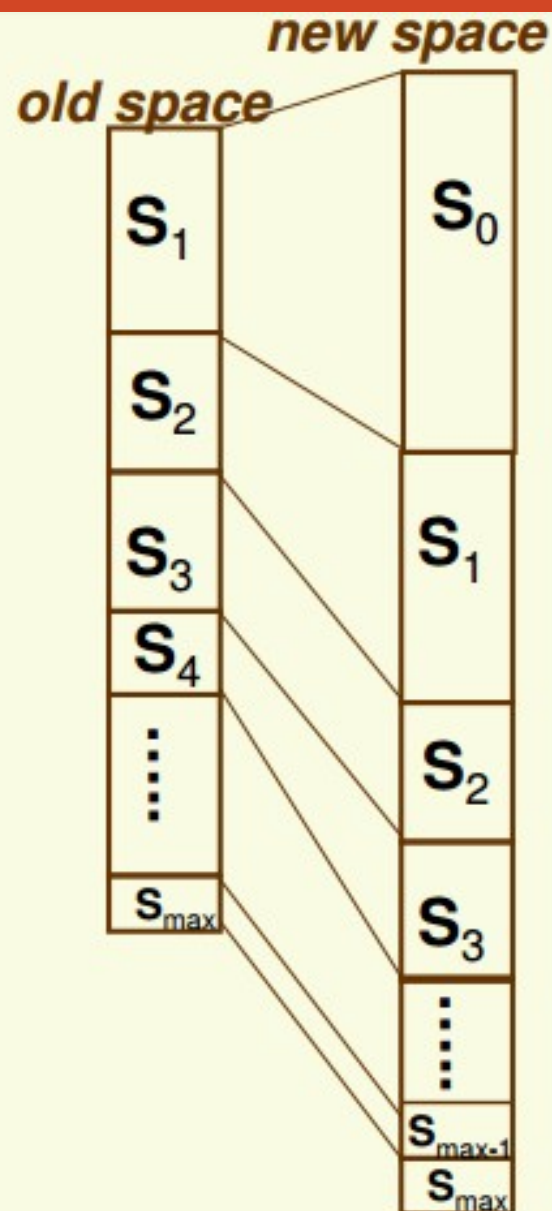
- Back to fishing: you have caught 10 Carp, 3 Cod, 2 tuna, 1 trout, 1 salmon, 1 eel; 20 species unseen
- How likely is it that next species is new? $3/18$
 - The probability of any individual unseen fish is $\frac{3}{18 \cdot 20}$
- What is the new probability of catching a trout?
 - Should be lower than $1/18^{\text{th}}$ to make room for unseen fish
 - Idea:
 - if we catch another trout, trout will occur with the rate of 2
 - According to our data, that is the probability of fish with rate 2 (occurring 2 times). Tuna occurs 2 times, so probability is $2/18$
 - Now spread the probability of $2/18$ over all species which occurred only once – 3 species
 - The probability of catching a fish which occurred 1 time already is $\frac{2}{18 \cdot 3}$

Smoothing: Good Turing

- In general, let r be the rate with which an n -gram occurs in the training data
 - Rate is the same thing as count
 - Example: if training data is {"a cow", "a train", "a cow", "do as", "to go", "let us", "to go"}, then the rate of "a cow" is 2 and the rate of "let us" is 1
- If an n -gram occurs with rate r , we used to get its probability as
 - r/N , where N is the size of the training data
 - We need to lower all the rates to make room for unseen n -grams
- In general, the number of n -grams which occur with rate $r+1$ is smaller than the number of grams which occur with rate r
- Idea: take the portion of probability space occupied by n -grams which occur with rate $r+1$ and divide it among the n -grams which occur with rate r

Smoothing: Good Turing

- Let S_r be the n-grams that occur r times in the training data
- Proportion of probability space occupied by n-grams in S_r in the new space = proportion of probability space occupied by n-grams in S_{r+1} in the new space
 - Spread evenly among all n-grams in S_r
- Note no space left for n-grams in S_{\max} , has to be fixed



Smoothing: Formula for Good Turing

- N_r be the number different n-grams that we saw in the training data exactly r times
 - Example: if training data is {"a cow", "a train", "a cow", "do as", "to go", "let us", "to go"}, then $N_1 = 3$ and $N_2 = 2$
 - In notation on previous slide, rN_r is the size of S_r
- Probability for any n-gram with rate r is estimated from the space occupied by n-grams with rate $r+1$
- Let N be the size of the training data. The probability space occupied by n-grams with rate $r+1$ is:

$$\frac{(r+1)N_{r+1}}{N}$$

- Spread this mass evenly among n-grams with rate r , there are N_r of them

$$\frac{(r+1)N_{r+1}}{N \cdot N_r}$$

- That is for a n-gram x that occurs r times, Good Turing estimate of probability is

$$P_{GT}(x) = (r+1) \frac{N_{r+1}}{N \cdot N_r}$$

Smoothing: Good Turing

$$P_{GT}(w_1 \dots w_n) = \frac{1}{N} \cdot \underbrace{\frac{(r+1)N_{r+1}}{N_r}}_{r^*}$$

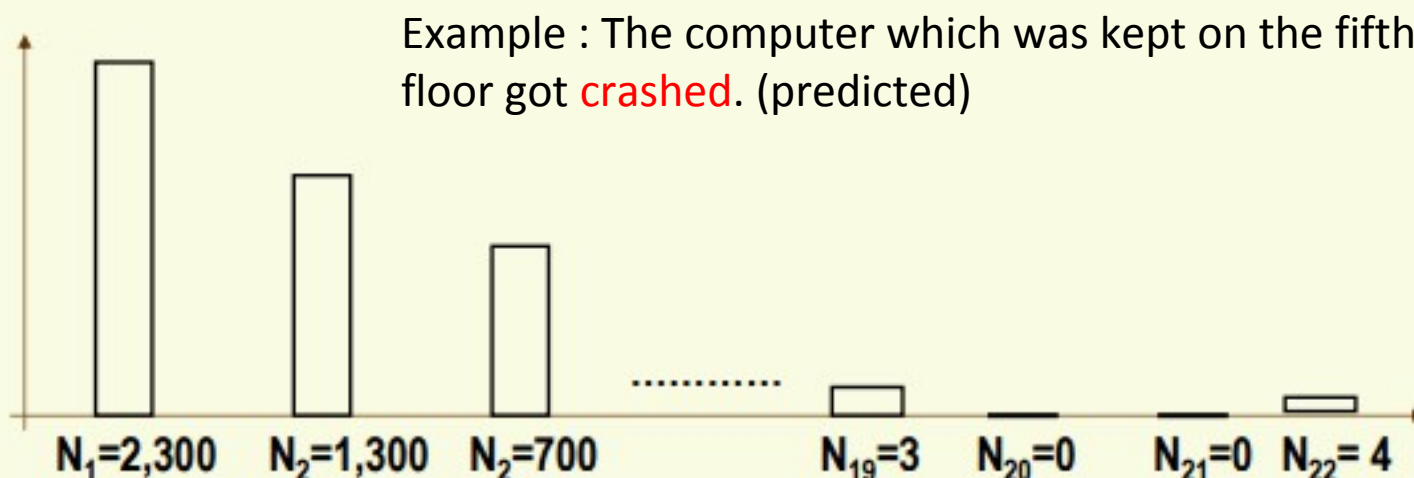
- Another way of looking at Good-Turing:

$$P_{MLE}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n)}{N} = \frac{r}{N}$$

- $P_{MLE}(w_1 \dots w_n) = 0$ for rate $r = 0$, need to increase it
 - at the expense of decreasing the rate of observed nGrams
 - if $r = 0$, new r^* should be larger
 - if $r \neq 0$, new r^* should be smaller
- This is exactly what Good-Turing does
 - For $r = 0$, $r^* = \frac{N_1}{N_0} > r$
 - For $r > 0$, $r^* = \frac{(r+1)N_{r+1}}{N_r}$
 - most likely $r^* < r$ since usually N_{r+1} is significantly less than N_r

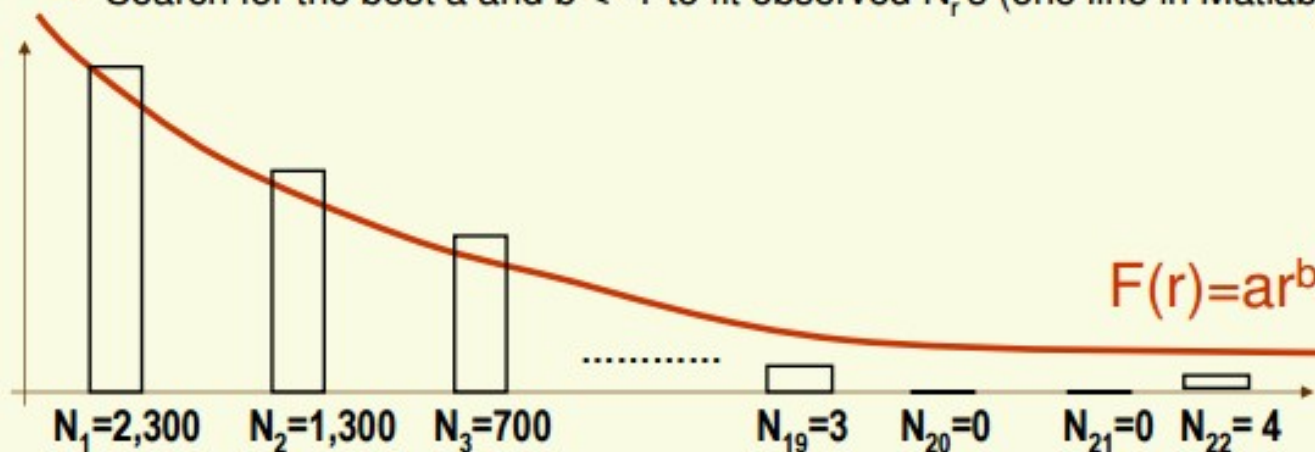
Smoothing: Fixing Good Turing

- That is for an n-gram x that occurs r times, Good Turing estimate of probability is
$$P_{GT}(x) = (r + 1) \frac{N_{r+1}}{N \cdot N_r}$$
- This works well except for high values of r
 - For high values of r , N_r is not reliable estimate of the number of n-grams that occur with rate r
 - In particular, for the most frequent r it completely fails since $N_{r+1}=0$
- The problem is that N_r is unreliable for high values of r



Smoothing: Fixing Good Turing

- The problem is that N_r is unreliable for high values of r
- Solution 1:
 - use P_{GT} for low values of r , say for $r < 10$
 - For n -grams with higher rates, use P_{MLE} which is reliable for higher values of r , that is $P_{MLE}(w_1 \dots w_n) = C(w_1 \dots w_n) / N$
- Solution 2:
 - Smooth out N_r 's by fitting a power law function $F(r) = ar^b$ (with $b < -1$) and use it when N_r becomes unreliable.
 - Search for the best a and $b < -1$ to fit observed N_r 's (one line in Matlab)



Good Turing vs. Add-One

$r = f_{\text{MLE}}$	$f_{\text{empirical}}$	f_{Lap}	f_{GT}
0	0.000027	0.000137	0.000027
1	0.448	0.000274	0.446
2	1.25	0.000411	1.26
3	2.24	0.000548	2.24
4	3.23	0.000685	3.24
5	4.21	0.000822	4.22
6	5.23	0.000959	5.19
7	6.21	0.00109	6.21
8	7.21	0.00123	7.24
9	8.26	0.00137	8.25

Smoothing: Fixing Good Turing

- Probabilities will not add up to 1, whether using Solution 1 or Solution 2 from the previous slide
- Have to renormalize all probabilities so that they add up to 1
 - Could renormalize all n-grams
 - Usually we renormalize only the n-grams with observed rates higher than 0
 - Suppose the total space for unseen n-grams is 1/20
 - renormalize the weight of the seen n-grams so that the total is 19/20

Question

NLTK

8)

Suppose you are reading an article on Natural Language Processing. Till now, you have read the words "language" - 8 times, "aspect" - 3 times, "processing" - 2 times, "extraction" - 2 times, "question" - once and "dialogue" - once. What are the Maximum Likelihood Estimate (MLE) probability ($P_{\text{processing}}$) and Good Turing probability ($P^*_{\text{GT}(\text{processing})}$) for reading "processing" as the next word?

1. 1/17, 2/17
2. 2/17, 1/17
3. 3/17, 2/17
4. 2/17, 1.5/17

- ☐ 1.
☐ 2.
☐ 3.
☐ 4.

No, the answer is incorrect.

Score: 0

Accepted Answers:

4.

9)

With the same setting as Question 8, calculate the MLE and Good Turing probabilities for reading "answering" as the next word:

1. 1/17, 2/17
2. 0, 1/17
3. 0, 2/17
4. 1/17, 1/17

- ☐ 1.
☐ 2.
☐ 3.
☐ 4.

No, the answer is incorrect.

Score: 0

Accepted Answers:

3.