

N-Grams and Language Models (Chapter 4)

Topics to be covered

- Word prediction
- N-grams
 - Counting and basic concepts
- Language Model
- Shannon's Method
 - Evaluation
- Smoothing for LM

Word Prediction

- Guess the next word...
 - *Please turn your homework ???*
 - *... I notice three guys standing on the ???*
- We can formalize this task using what are called **N-gram** models.

Word Prediction

- N-grams are **token sequences of length N** .
- *I notice three guys standing on the*
- 2-grams (aka bigrams)
 - (I notice), (notice three), (three guys), (guys standing), (standing on), (on the)
- What about 3-grams (aka **trigram**)
 - (I notice three), (notice three guys), (three guys standing), (guys standing on), (standing on the)

More examples

- Given knowledge of counts of N-grams such as these, we can guess likely next words in a sequence.
 - Predict the next word from the preceding N-1 words
- *I notice three guys standing on the*
 - Given a dataset (aka **corpus, corpora**), “one the street” 50 times, “on the computer” 2000
 - “Standing on the street” 40; “standing on the computer” 1

***N*-Gram Models**

- formally, we can use knowledge of the counts of *N*-grams to assess the **conditional probability** of candidate words as the next word in a sequence.
 - $P(\text{table} | \text{I notice three guys standing on the})$
- Or, we can use them to assess the **probability** of an entire sequence of words.
 - $P(\text{I notice three guys standing on the table})$
- Pretty much the same thing as we'll see...

Applications

- It turns out that being able to predict the next word (or any linguistic unit) in a sequence is an extremely useful thing to do.
- it lies at the core of the following applications
 - Automatic speech recognition
 - Handwriting and character recognition
 - Spelling correction
 - Machine translation
 - And many more.

Counting

- Simple counting lies at the core of any probabilistic approach. So let's first take a look at what we're counting.
 - *He stepped out into the hall, was delighted to encounter a water brother.*
 - 13 tokens, 15 if we include “,” and “.” as separate tokens.
 - Assuming we include the comma and period, how many bigrams are there?

Counting: Types and Tokens

- How about
 - *They picnicked by the pool, then lay back on the grass and looked at the stars.*
 - 18 tokens (again counting punctuation)
- But we might also note that “*the*” is used 3 times, so there are only 16 unique **types** (as opposed to **tokens**).
- In going forward, we’ll have occasion to focus on counting both types and tokens of both words and *N*-grams.

Counting: Corpora

- large collections of text
- Brown et al (1992) large corpus of English text
 - 583 million tokens
 - 293,181 types
- Google
 - Crawl of 1,024,908,267,229 English tokens
 - 13,588,391 types
 - That seems like a lot of types... After all, even large dictionaries of English have only around 500k types. Why so many here?
 - Numbers
 - Misspellings
 - Names
 - Acronyms
 - etc

Language Modeling

- $P(w_1, w_2 \dots w_{n-1} w_n)$, the probability of a sequence
- We can model the word prediction task as the conditional probability of a word given the previous words in the sequence
 - $P(w_n | w_1, w_2 \dots w_{n-1})$
- We'll call a statistical model that can assess these a *Language Model*
- How to compute
 - $P(\text{its water is so transparent that the})$
- Let's use the chain rule of probability

The Chain Rule

- Recall the definition of conditional probabilities

$$P(A | B) = \frac{P(A \wedge B)}{P(B)}$$

- Rewriting:

$$P(A \wedge B) = P(A | B)P(B)$$

- For sequences...

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

- In general

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

The Chain Rule

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2) \dots P(w_n|w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k|w_1^{k-1}) \end{aligned}$$

P(its water was so transparent)=

P(its)*

P(water|its)*

P(was|its water)*

P(so|its water was)*

P(transparent|its water was so)

Language Modeling

- How to calculate?
 - $P(\textit{the} \mid \textit{its water is so transparent that})$
- By definition of **conditional probabilities**
 $P(\textit{its water is so transparent that the})$
 $P(\textit{its water is so transparent that})$
We can get each of those from counts in a large corpus.

Very Easy Estimate

- How to estimate?
 - $P(\text{the } | \text{ its water is so transparent that})$

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

Very Easy Estimate

- According to Google those counts are 5/9.
 - that's **not terribly convincing** due to the small numbers involved.
- Unfortunately, for most sequences and for most text collections we won't get good estimates from this method.
 - What we're likely to get is 0. Or worse 0/0.

Summary

- There are still a lot of possible sentences
- In general, we'll never be able to get **enough data** to compute the statistics for **those longer** prefixes
 - $P(\text{the} \mid \text{its water is so transparent that})$
- Same problem we had for the strings themselves
 - $P(\text{its water is so transparent that } \textit{the})$

Independence Assumption

- Make the simplifying assumption
 - $P(\text{lizard} | \text{the, other, day, I, was, walking, along, and, saw, a}) = P(\text{lizard} | \text{a})$
- Or maybe
 - $P(\text{lizard} | \text{the, other, day, I, was, walking, along, and, saw, a}) = P(\text{lizard} | \text{saw, a})$
- Or maybe ??
- That is, the probability in question is independent of its earlier history.

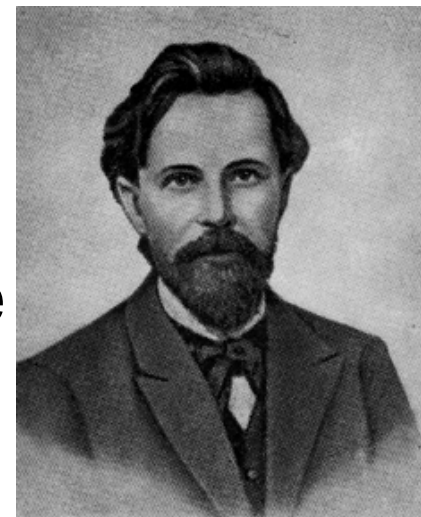
Independence Assumption

- This particular kind of independence assumption is called a *Markov assumption* after the Russian mathematician Andrei Markov.
- Markov Assumption: So for each component in the product replace with the approximation (assuming a prefix of N)

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

- Bigram version

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



Estimating Bigram Probabilities

- The Maximum Likelihood Estimate (MLE)

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

The maximum likelihood estimate of $p(w_i | w_{i-1})$ from a training set T (corpus)

An Example

Training data (corpus)

- $\langle s \rangle$ I am Sam $\langle /s \rangle$
- $\langle s \rangle$ Sam I am $\langle /s \rangle$
- $\langle s \rangle$ I do not like green eggs and ham $\langle /s \rangle$

$$\begin{aligned} P(I | \langle s \rangle) &= \frac{2}{3} = .67 & P(\text{Sam} | \langle s \rangle) &= \frac{1}{3} = .33 & P(\text{am} | I) &= \frac{2}{3} = .67 \\ P(\langle /s \rangle | \text{Sam}) &= \frac{1}{2} = 0.5 & P(\text{Sam} | \text{am}) &= \frac{1}{2} = .5 & P(\text{do} | I) &= \frac{1}{3} = .33 \end{aligned}$$

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

- But it is the **estimate** that makes it **most likely** that “I” will occur after $\langle s \rangle$ with a probability .67

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

- Out of 9222 sentences
 - Eg. "I want" occurred 827 times

	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

Bigram Probabilities

- $P(\text{want} \mid i)$
- Divide bigram counts by prefix unigram counts to get probabilities. $P(\text{want} \mid i)$

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

	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

- $P(<s> \text{ I want english food } </s>) =$
 $P(i|<s>)*$
 $P(\text{want}|I)*$
 $P(\text{english}|\text{want})*$
 $P(\text{food}|\text{english})*$
 $P(</s>|\text{food})*$
 $=.000031$

Kinds of Knowledge

- As crude as they are, N -gram probabilities capture a range of interesting facts about language.
- $P(\text{english}|\text{want}) = .0011$
- $P(\text{chinese}|\text{want}) = .0065$
- $P(\text{to}|\text{want}) = .66$
- $P(\text{eat} | \text{to}) = .28$
- $P(\text{food} | \text{to}) = 0$
- $P(\text{want} | \text{spend}) = 0$

World knowledge

Syntax

Shannon's Method

- Assigning probabilities to sentences is all well and good.
- A task is to turn the model around and use it to **generate** random sentences that are *like* the sentences from which the model was derived.
- Generally attributed to Claude Shannon.



Shannon's Method

- Sample a random bigram ($\langle s \rangle$, w) according to its probability
- Now sample a random bigram (w , x) according to its probability
 - Where the prefix w matches the suffix of the first.
- And so on until we randomly choose a (y , $\langle /s \rangle$)
- Then string the words together
- $\langle s \rangle$ I

I want

want to

to eat

eat Chinese

Chinese food

food $\langle /s \rangle$

Shakespeare

Unigram	<ul style="list-style-type: none"> • To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have • Every enter now severally so, let • Hill he late speaks; or! a more to leg less first you enter • Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like
Bigram	<ul style="list-style-type: none"> • What means, sir. I confess she? then all sorts, he is trim, captain. • Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow. • What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first gentleman? • Enter Menenius, if it so many good direction found'st thou art a strong upon command of fear not a liberal largess given away, Falstaff! Exeunt
Trigram	<ul style="list-style-type: none"> • Sweet prince, Falstaff shall die. Harry of Monmouth's grave. • This shall forbid it should be branded, if renown made it empty. • Indeed the duke; and had a very good friend. • Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.
Quadrigram	<ul style="list-style-type: none"> • King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in; • Will you not tell me who I am? • It cannot be but so. • Indeed the short and the long. Marry, 'tis a noble Lepidus.

Shakespeare as a Corpus

- $N=884,647$ tokens, $V=29,066$
- Shakespeare produced 300,000 bigram types out of $V^2= 844$ million possible bigrams...
 - So, 99.96% of the possible bigrams were never seen (have zero entries in the table)
- Quadrigrams are worse: What's coming out looks like Shakespeare because it ***is*** Shakespeare

The Wall Street Journal is Not Shakespeare

unigram: Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives

bigram: Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her

trigram: They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions

Questions

- Well Shannon's game gives us an intuition—higher order models look better
 - Quadrigrams can be used for word prediction? Can we make that notion operational?
 - The higher order models are likely to be pretty useless, especially from a small corpus
 - The probability $p(\text{Indeed the short and})$ is very small

Evaluation

- How do we know if our models are any good?
- Standard method
 - Train parameters of our model on a **training set**.
 - Look at the models performance on **some new data**
 - So use a **test set**. A dataset which is different than our training set
 - Then we need an **evaluation metric** to tell us how well our model is doing on the test set.
 - One such metric is **perplexity** (to be introduced below)

Perplexity

- Perplexity is the probability of the test set (assigned by the language model), normalized by the number of words:

$$\begin{aligned} \text{PP}(W) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}} \end{aligned}$$

N : number of words in a test data w_1, \dots, w_N is the test data

- Chain rule:

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

- For bigrams:

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

Lower perplexity means a better model

- Minimizing perplexity is the same as maximizing probability
 - **The best language model is one that best predicts an unseen test set**
- Given a training corpus to build LM, and test data. We can compute perplexity
- $PP(I \text{ want english food}) =$

$$\sqrt[N]{\frac{1}{P(\text{want}|I) * P(\text{english}|\text{want}) * P(\text{food}|\text{english})}}$$

Unknown Words– out of vocabulary (OOV)

- But once we start looking at test data, we'll run into words that we haven't seen before (pretty much regardless of how much training data you have.)
- With an *Open Vocabulary* task
 - Create an unknown word token <UNK>
 - Training of <UNK> probabilities
 - Create a fixed lexicon L , of size V
 - From a dictionary or
 - A subset of terms from the training set
 - At text normalization phase, any training word not in L changed to <UNK>
 - Now we count that like a normal word
 - At test time
 - Use UNK counts for any word not in training

Zero Counts

- Back to Shakespeare
 - Recall that Shakespeare produced 300,000 bigram types out of $V^2 = 844$ million possible bigrams...
 - So, 99.96% of the possible bigrams were never seen (have zero entries in the table)
 - Does that mean that any sentence that contains one of those bigrams should have a probability of 0?
- Some of those zeros are really zeros...
 - Things that really can't or shouldn't happen.
- On the other hand, some of them are just rare events.
 - If the training corpus had been a little bigger they would have had a count (probably a count of 1!).

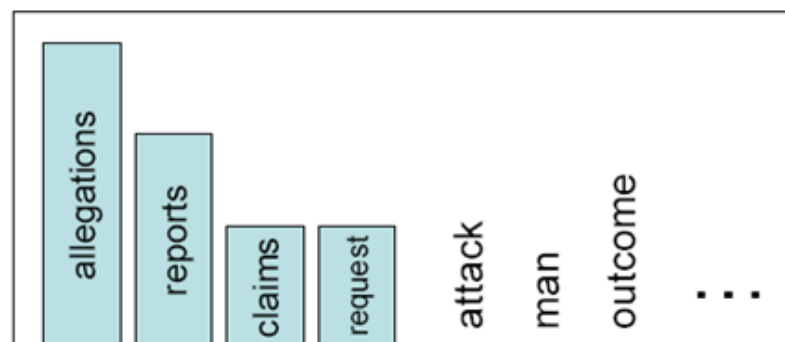
Zero Counts

- Zipf's Law (long tail phenomenon):
 - A small number of events occur with high frequency
 - A large number of events occur with low frequency
 - You can quickly collect statistics on the high frequency events
 - You might have to wait an arbitrarily long time to get valid statistics on low frequency events
- Result:
 - Our estimates are sparse! We have no counts at all for the vast bulk of things we want to estimate!
- Answer:
 - Estimate the likelihood of unseen (zero count) N-grams!

Smoothing is like Robin Hood: Steal from the rich and give to the poor (in probability mass)

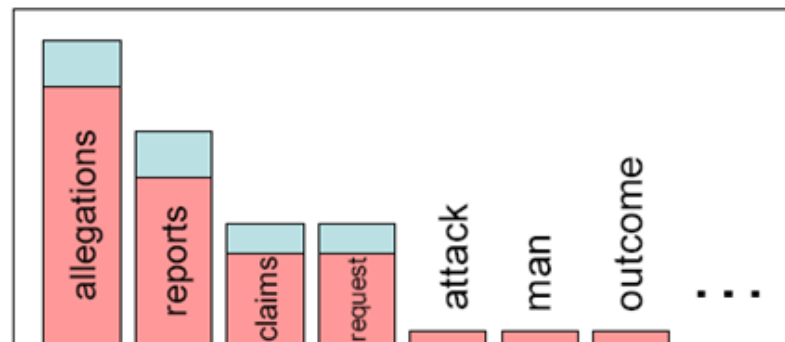
- We often want to make predictions from sparse statistics:

$P(w \mid \text{denied the})$
3 allegations
2 reports
1 claims
1 request
7 total



- Smoothing flattens spiky distributions so they generalize better

$P(w \mid \text{denied the})$
2.5 allegations
1.5 reports
0.5 claims
0.5 request
2 other
7 total



- Very important all over NLP, but easy to do badly!

Laplace Smoothing

- Also called add-one smoothing
- Just add one to all the counts!
- Very simple



- MLE estimate: $P(w_i) = \frac{c_i}{N}$

- Laplace estimate: $P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$

Laplace-Smoothed Bigram Counts

	i	want	to	eat	chinese	food	lunch	spend
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

Laplace-Smoothed Bigram Probabilities

$$P^*(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + V}$$

	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

Big Change to the Counts!

- $P(\text{to}|\text{want})$ from .66 to .26!
- Despite its flaws Laplace (add-k) is however still used to smooth other probabilistic models in NLP, especially
 - in domains where the number of zeros isn't so huge.

Better Smoothing

- Intuition used by many smoothing algorithms
 - Good-Turing
 - Kneser-Ney
 - Witten-Bell
- Is to use the count of things we've seen once to help estimate the count of things we've never seen



Backoff and Interpolation

- Another really useful source of knowledge
- If we are estimating:
 - trigram $p(z|x,y)$
 - but $\text{count}(xyz)$ is zero
- Use info from:
 - Bigram $p(z|y)$
- Or even:
 - Unigram $p(z)$
- How to combine this trigram, bigram, unigram info in a valid fashion?

Backoff Vs. Interpolation

- **Backoff:** use trigram if you have it, otherwise bigram, otherwise unigram
- **Interpolation:** mix all three
 - Simple interpolation

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

- More complicated interpolation (read book)

Practical Issues

- We do everything in log space
 - Avoid underflow
 - (also adding is faster than multiplying)
- $P(<s> \text{ I want english food } </s>)$
= $\frac{P(i|<s>)^* \quad P(\text{want}|I)^* \quad P(\text{english}|\text{want})^*}{P(\text{food}|\text{english})^* \quad P(</s>|\text{food})^*}$

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

Language Modeling Toolkits

Chapter
4.8 in
book

- SRILM
- CMU-Cambridge LM Toolkit
- These toolkits are publicly available
- Can use it to get N-gram models
- Lots of parameters (need to know the theory!)
- Standard N-gram format: ARPA language

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

CMU-Cambridge LM: <http://mi.eng.cam.ac.uk/~prc14/toolkit.html>

Google N-Gram Release

All Our N-gram are Belong to You

By Peter Norvig - 8/03/2006 11:26:00 AM

Posted by Alex Franz and Thorsten Brants, Google Machine Translation Team

Here at Google Research we have been using word [n-gram models](#) for a variety of R&D projects, such as [statistical machine translation](#), speech recognition, [spelling correction](#), entity detection, information extraction, and others. While such models have usually been estimated from training

to share this enormous dataset with everyone. We processed 1,024,908,267,229 words of running text and are publishing the counts for all 1,176,470,663 five-word sequences that appear at least 40 times. There are 13,588,391 unique words, after discarding words that appear less than 200 times.

<http://googleresearch.blogspot.sg/2006/08/all-our-n-gram-are-belong-to-you.html>

Google N-Gram Release

- 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
- serve as the indicator 120
- serve as the indicators 45
- serve as the indispensable 111
- serve as the indispensable 40
- serve as the individual 234

Google Caveat

- Remember the lesson about test sets and training sets... Test sets should be similar to the training set (drawn from the same distribution) for the probabilities to be meaningful.
- So... The Google corpus is fine if your application deals with arbitrary English text on the Web.
- If not then a smaller domain specific corpus is likely to yield better results.

Summary and sources

- Word prediction
 - N-grams
 - Counting and basic concepts
 - Language Model
 - Shannon's Method
 - Evaluation
 - Smoothing for LM
-
- **Great sources for NLP tools!**
 - <http://nlp.stanford.edu/links/statnlp.html>