

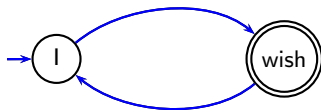
Language modelling

October 1, 2018

- Assign a probability to a sentence
- Applications
 - 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
 $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$
 - Text summarization, question-answering, ...

What is a language model?

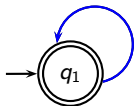
We can view a **finite state automaton** as a **deterministic** language model.



I wish I wish I wish I wish ...

- Cannot generate: “wish I wish” or “I wish I”
- Our basic model: each document was generated by a different automaton like this
- Except that these automata are probabilistic

A probabilistic language model



w	$P(w q_1)$	w	$P(w q_1)$
STOP	0.2	toad	0.01
the	0.2	said	0.03
a	0.1	likes	0.02
frog	0.01	that	0.04
	

- This is a one-state probabilistic finite-state automaton – a **unigram language model** – and the state emission distribution for its one state q_1 .
- STOP is not a word, but a special symbol indicating that the automaton stops.
- *frog said that toad likes frog STOP*
- $P(\text{string}) = 0.01 \times 0.03 \times 0.04 \times 0.01 \times 0.02 \times 0.01 \times 0.2 = 0.0000000000048$

A different language model for each document

language model of d_1				language model of d_2			
w	$P(w .)$	w	$P(w .)$	w	$P(w .)$	w	$P(w .)$
STOP	.2	toad	.01	STOP	.2	toad	.02
the	.2	said	.03	the	.15	said	.03
a	.1	likes	.02	a	.08	likes	.02
frog	.01	that	.04	frog	.01	that	.05
	

frog said that toad likes frog STOP

$$P(\text{query}|M_{d_1}) = 0.01 \times 0.03 \times 0.04 = 0.0000000000048 = 4.8 \times 10^{-12}$$

$$P(\text{query}|M_{d_2}) = 0.01 \times 0.02 \times 0.03 \times 0.05 = 0.0000000000030 = 3.0 \times 10^{-12}$$

$$P(\text{query}|M_{d_1}) < P(\text{query}|M_{d_2})$$

Thus, document d_2 is "more relevant" to the query "frog said that toad likes frog STOP" than d_1 is.

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

$$P(S) = P(w_1, w_2, w_3, w_4, w_5 w_n) \quad (1)$$

- Related task: probability of an upcoming word:

$$P(w_5 | w_1, w_2, w_3, w_4) \quad (2)$$

- A model that computes either of these: $P(S)$ or $P(w_n | w_1, w_2 w_{n-1})$ is called a language model.
- Better: the grammar
- But language model or LM is standard

How to compute $P(S)$

- How to compute this joint probability:

$$P(\textit{its, water, is, so, transparent, that})$$

- Intuition: lets rely on the Chain Rule of Probability

The chain rule

- Conditional probabilities

$$P(A, B) = P(A \text{ and } B) = P(A)P(B|A) = P(B)P(A|B) \quad (3)$$

- More variables:

$$P(A, B, C, D) = P(A)P(B|A)P(C|A, B)P(D|A, B, C) \quad (4)$$

- The Chain Rule in General

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

- For a sentence:

$$P(w_1 w_2 \dots w_n) = \prod_{i=1}^n P(w_i | w_1 w_2 \dots w_{i-1})$$

-

$$\begin{aligned} P(\textit{its water is so transparent}) &= P(\textit{its}) \\ &\times P(\textit{water}|\textit{its}) \\ &\times P(\textit{is}|\textit{its water}) \\ &\times P(\textit{so}|\textit{its water is}) \\ &\times P(\textit{transparent}|\textit{its water is so}) \end{aligned} \quad (5)$$

How to estimate these probabilities

- Count and divide?

$$P(\textit{the}|\textit{its water is so transparent that}) \quad (6)$$

$$= \frac{\textit{count}(\textit{its water is so transparent that the})}{\textit{count}(\textit{its water is so transparent that})} \quad (7)$$

- Too many possible sentences
- Not enough data for estimating

What is language model? predict the probability of a sequence of words.
n-gram language model:

$$p(w_1, w_2, \dots, w_T) = \prod_i p(w_i | w_{i-1}, \dots, w_{i-n+1}) \quad (8)$$

Derived using chain rule and Markov assumption.

$$p(w_t | w_{t-n+1}, \dots, w_{t-1}) = \frac{\text{Count}(w_{t-n+1}, \dots, w_t)}{\text{Count}(w_{t-n+1}, \dots, w_{t-1})} \quad (9)$$

The problem of this approach: scarcity of data

Markov assumption

- Simplifying assumption:

$$P(\textit{the}|\textit{its water is so transparent that}) \approx P(\textit{the}|\textit{that})$$

- Or maybe

$$P(\textit{the}|\textit{its water is so transparent that}) \approx P(\textit{the}|\textit{transparent that})$$

- Markov assumption

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i | w_{i-k} \dots w_{i-1})$$

In other words, we approximate each component in the product

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx \prod_i P(w_i | w_{i-k} \dots w_{i-1})$$

Simplest case: the unigram model

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

Some automatically generated sentences from a unigram model:

- fifth, an, of, futures, the, an, incorporated, a, a, the, inflation, most, dollars, quarter, in, is, mass
- thrift, did, eighty, said, hard, 'm, july, bullish
- that, or, limited, the

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

Some automatically generated sentences from a bigram model:

texaco rose one in this issue is pursuing growth
in a boiler house said mr. gurria mexico 's motion
control proposal without permission from five
hundred fifty five yen

outside new car parking lot of the agreement
reached

this would be a record november

- We can extend to trigrams, 4-grams, 5-grams
- In general this is an insufficient model of language
- because language has long-distance dependencies:
*The **computer** which I had just put into the machine room on the fifth floor **crashed**.*
- But we can often get away with N-gram models

The Maximum Likelihood Estimate

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

- $\langle s \rangle$ *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* $\langle /s \rangle$

$$P(I|\langle s \rangle) = \frac{2}{3} \quad P(am|I) = \frac{2}{3} \quad P(Sam|am) = \frac{1}{2} \quad P(\langle /s \rangle | Sam) = \frac{1}{2}$$

More examples: Berkeley Restaurant Project sentences

- Example sentences

can you tell me about any good cantonese restaurants
close by

mid priced thai food is what im looking for

tell me about chez panisse

can you give me a listing of the kinds of food that are avail

im looking for a good place to eat breakfast

when is caffe venezia open during the day

Raw bigram counts

Out of 9222 sentences

	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

Normalized bigram

unigram:

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

normalized by the unigram:

	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

e.g., $i \text{ want} = 827 / 2533 = 0.33$

Bigram estimates of sentence probability

$$\begin{aligned} &P(\langle s \rangle I \text{ want english food } \langle /s \rangle) \\ &= P(I | \langle s \rangle) \\ &\times P(\text{want} | I) \\ &\times P(\text{english} | \text{want}) \\ &\times P(\text{food} | \text{english}) \\ &\times P(\langle /s \rangle | \text{food}) \\ &= .000031 \end{aligned}$$

$$\begin{aligned} P(j | \langle s \rangle) &= .25 \\ P(\text{english} | \text{want}) &= .0011 \\ P(\text{chinese} | \text{want}) &= .0065 \\ P(\text{to} | \text{want}) &= .66 \\ P(\text{eat} | \text{to}) &= .28 \\ P(\text{food} | \text{to}) &= 0 \\ P(\text{want} | \text{spend}) &= 0 \end{aligned}$$

We do everything in log space

- Avoid underflow
- also adding is faster than multiplying

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

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 book ngram: <http://ngrams.googlelabs.com/>
<http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html>

The Shannon Game:

- How well can we predict the next word?
 - I always order pizza with cheese and mushrooms 0.1

pepperoni 0.1

anchovies 0.01

...

fried rice 0.0001

...

and 1e-100

The 33rd President of the US was ____

I saw a ____

- Unigrams won't work for this task.
- A better model of a text is one which assigns a higher probability to the word that actually occurs

Perplexity

- The best language model is one that best predicts an unseen test set (Gives the highest $P(\text{sentence})$)
- Perplexity is the inverse probability of the test set, normalized by the number of words:

$$PP(S) = P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} = \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}}$$

- Chain rule:

$$PP(S) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 w_2 \dots w_{i-1})}}$$

- For bigrams:

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

The Shannon Game intuition for perplexity

- How hard is the task of recognizing digits 0,1,2,3,4,5,6,7,8,9
 - Lets suppose a sentence consisting of random digits
 - What is the perplexity of this sentence according to a model that assign $P=1/10$ to each digit?

$$\begin{aligned} PP(S) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \left[\left(\frac{1}{10} \right)^N \right]^{-\frac{1}{N}} \\ &= \left(\frac{1}{10} \right)^{-1} \\ &= 10 \end{aligned} \tag{10}$$

- How hard is recognizing (30,000) names at Microsoft.
 - Perplexity = 30,000
- If a system has to recognize
 - Operator (1 in 4)
 - Sales (1 in 4)
 - Technical Support (1 in 4)
 - 30,000 names (1 in 120,000 each)
 - Perplexity is 53
- Perplexity is weighted equivalent branching factor

Lower perplexity = better model

- Training 38 million words, test 1.5 million words, WSJ

	Unigram	Bigram	Trigram
Perplexity	962	170	109

Unigram

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

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?

Trigram

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

King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'
Will you not tell me who I am?
It cannot be but so.
Indeed the short and the long. Marry, 'tis a noble Lepidus.

- $N=884,647$ 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 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

Using language models (LMs) for IR

- 1 LM = language model
- 2 We view the document as a generative model that generates the query.
- 3 What we need to do:
 - Define the precise generative model we want to use
 - Estimate parameters (different parameters for each document's model)
 - Smooth to avoid zeros
 - Apply to query and find document most likely to have generated the query
 - Present most likely document(s) to user
- 4 Note that 3 is similar to what we did in Naive Bayes.

- Training set:
 - denied the allegations
 - denied the reports
 - denied the claims
 - denied the request
- Test set
 - denied the offer
 - denied the loan
- $P(\text{offer}|\text{denied the}) = 0$

- Key intuition: A nonoccurring term is possible (even though it didn't occur), ...
- ... but no more likely than would be expected by chance in the collection.
- Notation: M_c : the collection model; cf_t : the number of occurrences of t in the collection; $T = \sum_t cf_t$: the total number of tokens in the collection.

-

$$\hat{P}(t|M_c) = \frac{cf_t}{T}$$

- Use $\hat{P}(t|M_c)$ to “smooth” $P(t|d)$ away from zero.