

n-gram Language Modeling

n-gram, rnn, transformer model

Introduction

Predicting Words

- The water of Walden Pond is beautifully ...
- blue
- green
- clear
- refrigerator
- that

Language Models

Systems that can predict upcoming words

- Can assign a probability to each potential next word
- Can assign a probability to a whole sentence

Why word prediction?

It's a helpful part of language tasks

- Grammar or spell checking
- Their are two midterms ~~Their~~ There are two midterms
- Everything has improve Everything has ~~improve~~ improved
- Speech recognition
- I will be back soonish I will be bassoon dish

Why word prediction?

It's how large language models (LLMs) work!

LLMs are trained to predict words

- **Left-to-right (autoregressive) LMs learn to predict next word**

LLMs generate text by predicting words

- **By predicting the next word over and over again**

Language Model (LM) more formally

Goal – compute the probability of a sentence or sequence of words W :

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

Related task – probability of an upcoming word:

$$P(w_5 | w_1, w_2, w_3, w_4) \quad \text{or} \quad P(w_n | w_1, w_2, \dots, w_{n-1})$$

An LM computes either of these:

$$P(W) \quad \text{or} \quad P(w_n | w_1, w_2, \dots, w_{n-1})$$

How to estimate these probabilities

Could we just count and divide?

$$P(\textit{blue}|\textit{The water of Wanden Pond is so beautifully}) = \frac{C(\textit{The water of Walden Pond is so beautifully blue})}{C(\textit{The water of Walden Pond is so beautifully})}$$

No! Too many possible sentences!

We'll never see enough data for estimating these

How to compute $P(W)$ or $P(w_n|w_1, \dots, w_{n-1})$

How to compute the joint probability $P(W)$:

$P(\text{The, water, of, Walden, Pond, is, so, beautifully, blue})$

Intuition: let's rely on the Chain Rule of Probability

The Chain Rule

Recall the definition of conditional probabilities

$$P(B|A) = P(A, B)/P(A) \quad \text{Rewriting:} \quad P(A, B) = P(A)P(B|A)$$

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

$$P(X_1, X_2, X_3, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots P(X_n|X_1, \dots, X_{n-1})$$

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

$$P(w_{1:n}) = P(w_1)P(w_2|w_1)P(w_3|w_{1:2}) \dots P(w_n | w_{1:n-1})$$

$$= \prod_{k=1}^n P(w_k | w_{1:k-1})$$

P("The water of Walden Pond") =

**P(The) × P(Water | The) × P(of | The Water) × P(Walden | The water of) ×
P(Pond | The water of Walden)**

Markov Assumption

Simplifying assumption:

$P(\text{blue} \mid \text{The Water of Walden Pond is so beautifully})$
 $\approx P(\text{blue} \mid \text{beautifully})$

$$P(w_n \mid w_{1:n-1}) \approx P(w_n \mid w_{n-1})$$

Bigram Markov Assumption

$$P(w_{1:n}) \approx \prod_{k=1}^n P(w_k | w_{k-1})$$

Instead of $\prod_{k=1}^n P(w_k | w_{1:k-1})$

More generally in N-gram, we approximate each component in the product

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

Problems with N-gram models

- N-grams can't handle long-distance dependencies:

“The soups that I made from that new cookbook I bought yesterday were amazingly delicious.”

- N-grams don't do well at modeling new sequences with similar meanings

The solution: Large language models

- can handle much longer contexts
- because of using embedding spaces, can model synonymy better, and generate better novel strings

Why N-gram models?

n-gram 모델은 오늘날 Transformer Network 기반의 LLM에 밀렸지만, 그 단순함과 실용성 덕분에 "해석 가능한 통계 기반 NLP"의 대표 도구로 여전히 유효함

카테고리	대표 사용 사례
언어 처리	자동 완성, 음성 인식, 스펠링 교정
검색 · 색인	character n-gram 색인 (언어 불문), 쿼리 확장
평가 지표	BLEU, ROUGE, ChrF 등
특성 추출	BoW / TF-IDF 기반 분류기 전처리
생물학 · 응용	DNA 서열 분석, 키보드 예측, 단순 코드 완성

Estimating N-gram Probabilities

2

Estimating bigram probabilities

The Maximum Likelihood Estimate

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

- w_{n-1} 로 시작하는 모든 pair 수에 대한 $w_{n-1}w_n$ pair 수의 비율

위 식은 다음과 같이 간단히 정리될 수 있음

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

An example

<s> I am Sam </s>

<s> Sam I am </s>

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

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

$$P(I|<s>) = \frac{2}{3} = 0.67 \quad P(Sam|<s>) = \frac{1}{3} = 0.33 \quad P(am|I) = \frac{2}{3} = 0.67$$

$$P(<s>|Sam) = \frac{1}{2} = 0.5 \quad P(Sam|am) = \frac{1}{2} = 0.5 \quad P(do|I) = \frac{1}{3} = 0.33$$

More examples: Berkeley Restaurant Project sentences

can you tell me about any good cantonese restaurants close by

tell me about chez panisse

i'm 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

Raw bigram probabilities

Normalize by unigrams:

I	want	to	eat	Chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

Results:

	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 estimate of sentence probabilities

$P(\langle s \rangle \text{ I want english food } \langle /s \rangle) =$

$P(\text{I} | \langle s \rangle)$

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

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

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

$\times P(\langle /s \rangle | \text{food})$

$= .000031$

What kinds of knowledge do N-grams represent?

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

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

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

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

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

$$P(\text{want} \mid \text{spend}) = 0$$

$$P(i \mid \langle s \rangle) = .25$$

Dealing with scale in large n-grams

LM probabilities are stored and computed in log format, i.e. log probabilities

This avoids underflow from multiplying many small numbers

$$\log(P_1 \times P_2 \times P_3 \times P_4) = \log P_1 + \log P_2 + \log P_3 + \log P_4$$

If we need probabilities, we can do one exp at the end

$$P_1 \times P_2 \times P_3 \times P_4 = \exp(\log P_1 + \log P_2 + \log P_3 + \log P_4)$$

Larger n-grams

- 4-grams, 5-grams
- Large datasets of large n-grams have been released
 - N-grams from Corpus of Contemporary American English (COCA) 1 billion words (Davies 2020)
 - Google Web 5-grams (Franz and Brants 2006) 1 trillion words)
 - Efficiency: quantize probabilities to 4-8 bits instead of 8-byte float
- Newest model: infini-grams (∞ -grams) (Liu et al 2024)
 - No precomputing! Instead, store 5 trillion words of web text in suffix arrays. Can compute n-gram probabilities with any n!

N-gram LM Toolkits

SRILM

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

KenLM

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

NLTK

Evaluation and Perplexity



How to evaluate N-gram models – Extrinsic (in-vivo) Evaluation

- To compare models A and B

1. Put each model in a real task

- Machine Translation, speech recognition, etc.

2. Run the task, get a score for A and for B

- How many words translated correctly

3. How many words transcribed correctly

- Compare accuracy for A and B

Intrinsic (in-vitro) evaluation

Extrinsic evaluation not always possible

- Expensive, time-consuming
- Doesn't always generalize to other applications

Intrinsic evaluation: **perplexity**

- Directly measures language model performance at predicting words.
- Doesn't necessarily correspond with real application performance
- But gives us a single general metric for language models
- **Useful for large language models (LLMs) as well as n-grams**

Training sets and test sets

- We train parameters of our model on a **training set**.
- We test the model's performance on data we haven't seen.
 - A **test set** is an unseen dataset; different from training set.
 - Intuition: we want to measure generalization to unseen data
 - An **evaluation metric** (like **perplexity**) tells us how well our model does on the test set.

Choosing training and test sets

- **If we're building an LM for a specific task**
 - The test set should reflect the task language we want to use the model for
- **If we're building a general-purpose model**
 - We'll need lots of different kinds of training data
 - We don't want the training set or the test set to be just from one domain or author or language.

Training on the test set

We can't allow test sentences into the training set

- **Or else the LM will assign that sentence an artificially high probability when we see it in the test set**
- **And hence assign the whole test set a falsely high probability.**
- **Making the LM look better than it really is**

This is called “Training on the test set”

Bad science!

Dev sets (Validation sets)

- If we test on the test set many times we might implicitly tune to its characteristics
 - Noticing which changes make the model better.
- So, we run on the test set only once, or a few times
- That means we need a third dataset:
 - A development test set or, devset.
 - We test our LM on the devset until the very end
 - And then test our LM on the test set once

Intuition of perplexity as evaluation metric 1

How good is our language model?

Intuition: A good LM prefers "real" sentences

- Assign higher probability to "real" or "frequently observed" sentences
- Assigns lower probability to "word salad" or "rarely observed" sentences?

Intuition of perplexity 2

Predicting upcoming words

- The Shannon Game: How well can we predict the next word?
- Once upon a ____
- That is a picture of a ____
- For breakfast I ate my usual ____

time	0.9
dream	0.03
midnight	0.02
...	
and	1e-100

- Unigrams are terrible at this game (Why?)
- A good LM is one that assigns a higher probability to the next word that actually occurs

Intuition of perplexity 3

- The best language model is one that best predicts the entire unseen test set
- We said: a good LM is one that assigns a higher probability to the next word that actually occurs.
- Let's generalize to all the words!
 - The best LM assigns high probability to the entire test set.
- When comparing two LMs, A and B
 - We compute $P_A(\text{test set})$ and $P_B(\text{test set})$
 - The better LM will give a higher probability to (=be less surprised by) the test set than the other LM.

Intuition of perplexity 4

Use perplexity instead of raw probability

- Probability depends on size of test set
 - Probability gets smaller the longer the text
 - Better: a metric that is per-word, normalized by length
- Perplexity is the inverse probability of the test set, normalized by the number of words

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

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

Intuition of perplexity

The inverse

- Perplexity is the inverse probability of the test set, normalized by the number of words

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

- (The inverse comes from the original definition of perplexity from cross-entropy rate in information theory)
- Probability range is $[0,1]$, perplexity range is $[1,\infty]$
- Minimizing perplexity is the same as maximizing probability

Intuition of perplexity 6 – N-grams

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

Chain rule:

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

Bigrams:

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

Intuition of perplexity 7 – Weighted average branching factor

Perplexity is also the weighted average branching factor of a language.

Branching factor: number of possible next words that can follow any word

Example: Deterministic language $L = \{\text{red, blue, green}\}$

Branching factor = 3 (any word can be followed by red, blue, green)

Now **assume LM A** where each word follows any other word with equal probability $\frac{1}{3}$

Given a test set $T = \text{"red red red red blue"}$

$$\text{Perplexity}_A(T) = P_A(\text{red red red red blue})^{-\frac{1}{5}} = \left(\left(\frac{1}{3} \right)^5 \right)^{-1/5} = \left(\frac{1}{3} \right)^{-1} = 3$$

- But now suppose red was very likely in training set, such that **for LM B**:
 - $P(\text{red}) = .8$ $p(\text{green}) = .1$ $p(\text{blue}) = .1$
- We would expect the probability to be higher, and hence the perplexity to be smaller:

$$\text{Perplexity}_B(T) = P_B(\text{red red red red blue})^{-1/5} = (.8 \times .8 \times .8 \times .8 \times .1)^{-1/5} = 1.89$$

Holding test set constant: Lower perplexity = better LM

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

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

Sampling and Generation



The Shannon (1948) Visualization Method Sample words from an LM

Unigram:

REPRESENTING AND SPEEDILY IS AN GOOD APT OR COME CAN DIFFERENT NATURAL HERE
HE THE A IN CAME THE TO OF TO EXPERT GRAY COME TO FURNISHES THE LINE MESSAGE
HAD BE THESE.

Bigram:

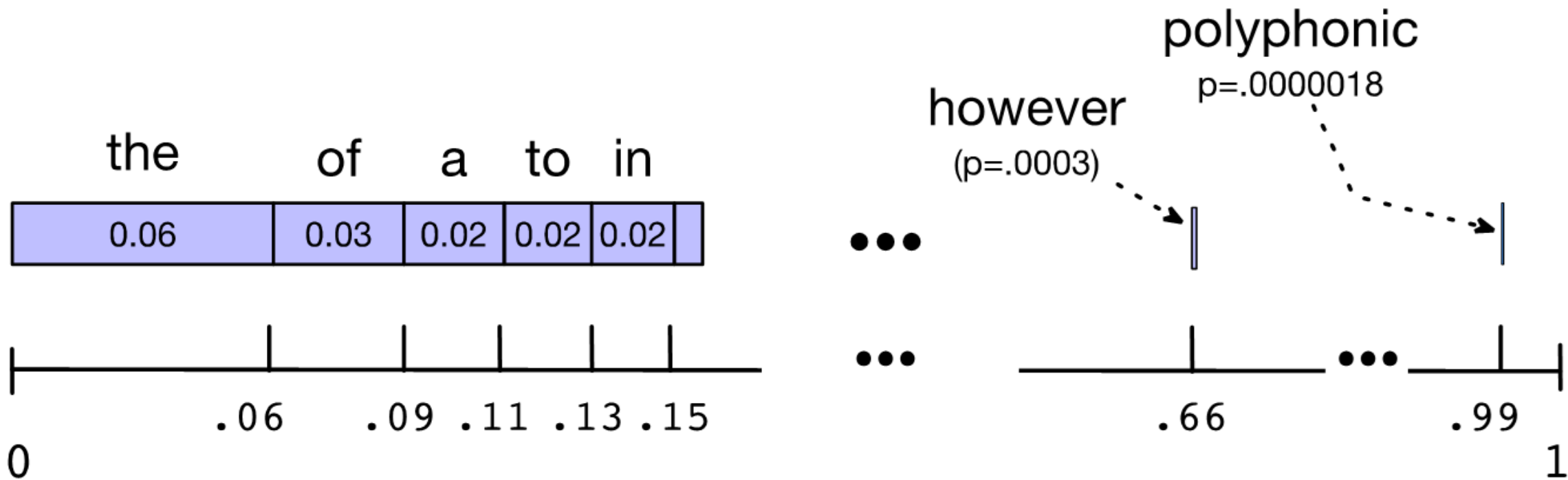
THE HEAD AND IN FRONTAL ATTACK ON AN ENGLISH WRITER THAT THE CHARACTER OF
THIS POINT IS THEREFORE ANOTHER METHOD FOR THE LETTERS THAT THE TIME OF WHO
EVER TOLD THE PROBLEM FOR AN UNEXPECTED.

How Shannon sampled those words in 1948



"Open a book at random and select a letter at random on the page. This letter is recorded. The book is then opened to another page, and one reads until this letter is encountered. The succeeding letter is then recorded. Turning to another page this second letter is searched for and the succeeding letter recorded, etc."

Sampling a word from a distribution



Visualizing Bigrams the Shannon Way

Choose a random bigram ($\langle s \rangle$, w)

according to its probability $p(w|\langle s \rangle)$

Now choose a random bigram (w , x) according

to its probability $p(x|w)$

And so, on until we choose $\langle /s \rangle$

Then string the words together

$\langle s \rangle$ I
 I want
 want to
 to eat
 eat Chinese
 Chinese food
 food $\langle /s \rangle$
 I want to eat Chinese food

Note: there are other sampling methods

Used for neural language models

Many of them avoid generating words from the very unlikely tail of the distribution

We'll discuss when we get to neural LM decoding:

- Temperature sampling
- Top-k sampling
- Top-p sampling

Approximating Shakespeare

1 gram	<p>–To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have</p> <p>–Hill he late speaks; or! a more to leg less first you enter</p>
2 gram	<p>–Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.</p> <p>–What means, sir. I confess she? then all sorts, he is trim, captain.</p>
3 gram	<p>–Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.</p> <p>–This shall forbid it should be branded, if renown made it empty.</p>
4 gram	<p>–King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;</p> <p>–It cannot be but so.</p>

This figure shows some random sentences generated from unigram, bigram, trigram, and 4- gram models trained on the works of William Shakespeare.

The longer the context, the more coherent the sentences. In these unigram sentences, there's no coherent relation between words. The bigram sentences have some local word-to-word coherence, and sensible punctuation. The 3 and 4-gram sentences are beginning to look a lot like Shakespeare. In fact, they look a little too much like Shakespeare. The words It cannot be but so are directly from King John.

Shakespeare as 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)**
- **That sparsity is even worse for 4-grams, explaining why our sampling generated actual Shakespeare.**

Just choosing one author, Shakespeare, makes for a pretty small corpus; under a million words from a vocabulary of 29,000 words. The vocab 29,000 squared is somewhat under a billion. So that means there are a billion possible bigrams whose counts have to come from only a million word corpus, and 10^{17} possible 4-grams. So, when training on any corpus, the vast majority of all the n-gram counts will be zero!

The Wall Street Journal is not Shakespeare

1
gram

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

2
gram

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

3
gram

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

To get an idea of the dependence of a grammar on its training set, let's look at an n-gram grammar trained on a completely different corpus: the Wall Street Journal (WSJ) newspaper. Shakespeare and the Wall Street Journal are both English, so we might expect some overlap between our n-grams. There is no overlap even in small phrases, let alone entire sentences. Statistical models are use- less as predictors when the training sets and the test sets are as different as Shakespeare and WSJ.

Can you guess the author?

- These 3-gram sentences are sampled from an LM trained on who?
1. 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 gram Brazil on market conditions
 2. This shall forbid it should be branded, if renown made it empty.
 3. "You are uniformly charming!" cried he, with a smile of associating and now and then I bowed and they perceived a chaise and four to wish for.

And just to make this absolutely clear, look at these 3 sentences. It should be very clear which one is sampled from the WSJ, which from Shakespeare, and which from Jane Austen.

Choosing training data

If task-specific, use a training corpus that has a similar genre to your task.

- **If legal or medical, need lots of special-purpose documents**

Make sure to cover different kinds of dialects and speaker/authors.

- **Example: African-American Vernacular English (AAVE)**
 - **One of many varieties that can be used by African Americans and others**
 - **Can include the auxiliary verb finna that marks immediate future tense:**
 - **"My phone finna die"**

To deal with this problem we should always use a training corpus that has a similar genre to whatever task we are trying to accomplish. To build a language model for translating legal documents, we need a training corpus of legal documents. To build a language model for a question-answering system, we need a training corpus of questions. It is equally important to get training data in the appropriate dialect or variety, especially when processing social media posts or spoken transcripts. For example some tweets will use features of African American English (AAE)—the name for the many varieties of language used in African American communities. Such features can include words like finna—an auxiliary verb that marks immediate future tense.

The perils of overfitting

N-grams only work well for word prediction if the test corpus looks like the training corpus

- But even when we try to pick a good training corpus, the test set will surprise us!
- We need to train robust models that generalize!

One kind of generalization: Zeros

- Things that don't ever occur in the training set
 - But occur in the test set
- In summary, n-grams only work well if the training set is chosen carefully to provide the right kind of language for the test setting. But matching genres and dialects is still not sufficient. Our models may still be subject to the problem of sparsity! Even a good training set might simply be missing some n-grams that appear in the test set.

Zeros

- Training set:
 - ... ate lunch
 - ... ate dinner
 - ... ate a
 - ... ate the

$$P(\text{"breakfast"} \mid \text{ate}) = 0$$

- Test set
 - ... ate lunch
 - ... ate breakfast

For example, consider a training set that has the phrases "ate lunch", "ate dinner", "ate a", and "ate the", but happens not to have had the phrase "ate breakfast". The $P(\text{"breakfas"} \mid \text{ate})$ in the test set will be 0!

Zero probabilities bigrams

Bigrams with zero probability

- Will hurt our performance for texts where those words appear!
- And mean that we will assign 0 probability to the test set!

And hence we cannot compute perplexity (can't divide by 0)!

These zeros are a problem for two reasons.

- **First**, they mean we are underestimating the probability of all sorts of words that might occur, which will hurt the performance of any application we want to run on this data.
- **Second**, if the probability of any word in the test set is 0, the entire probability of the test set is 0. And perplexity is based on the inverse probability of the test set. Thus, if some words have zero probability, we can't compute perplexity at all, since we can't divide by 0! The solution to this problem is called smoothing, and we'll see it in the next lecture!

Smoothing, Interpolation, and Backoff

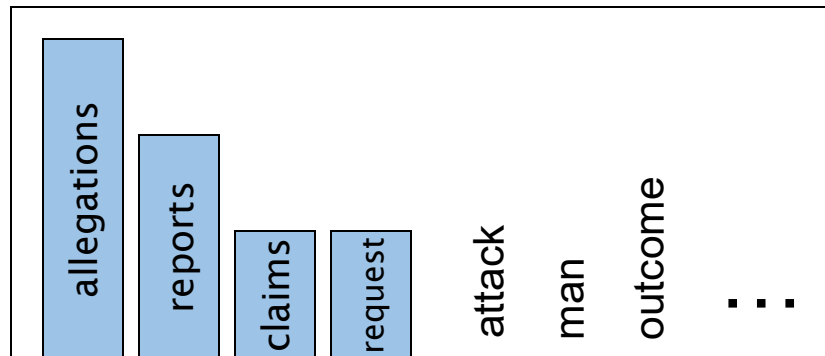


The intuition of smoothing (from Dan Klein)

When we have sparse statistics:

$P(w \mid \text{denied the})$

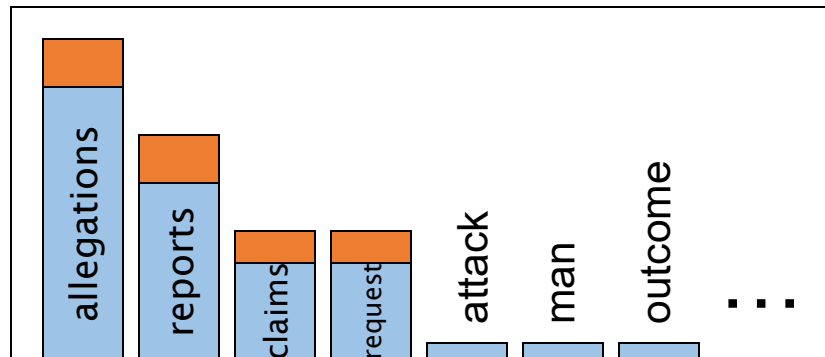
- 3 allegations
- 2 reports
- 1 claims
- 1 request
- 7 total



Steal probability mass to generalize better

$P(w \mid \text{denied the})$

- 2.5 allegations
- 1.5 reports
- 0.5 claims
- 0.5 request
- 2 other
- 7 total



Add-one estimation

Also called **Laplace smoothing**

Pretend we saw each word one more time than we did

Just add one to all the counts!

- MLE estimate:

$$P_{MLE}(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

- Add-1 estimate:

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

Maximum Likelihood Estimates

The maximum likelihood estimate

- of some parameter of a model M from a training set T
- maximizes the likelihood of the training set T given the model M

Suppose the word “bagel” occurs 400 times in a corpus of a million words

What is the probability that a random word from some other text will be “bagel”?

MLE estimate is $400/1,000,000 = .0004$

This may be a bad estimate for some other corpus

- But it is the estimate that makes it most likely that “bagel” will occur 400 times in a million word corpus.

Berkeley Restaurant Corpus: 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

	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

Laplace-smoothed bigrams

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

	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

Reconstructed counts

$$c^*(w_{n-1}w_n) = \frac{[C(w_{n-1}w_n) + 1] \times C(w_{n-1})}{C(w_{n-1}) + V}$$

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

Compare with raw bigram counts

	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

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

Add-1 estimation is a blunt instrument

So, add-1 isn't used for N-grams:

- Generally, we use interpolation or backoff instead

But add-1 is used to smooth other NLP models

- For text classification
- In domains where the number of zeros isn't so huge.

Backoff and Interpolation

Sometimes it helps to use less context

- Condition on less context for contexts you know less about

Backoff:

- use trigram if you have good evidence,
- otherwise, bigram, otherwise unigram

Interpolation:

- mix unigram, bigram, trigram

Interpolation works better

Linear Interpolation

- Simple interpolation

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

$$\sum_i \lambda_i = 1$$

- Lambdas conditional on context:

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

How to set λ s for interpolation?

Use a held-out corpus



Choose λ s to maximize probability of held-out data:

- Fix the N-gram probabilities (on the training data)
- Then search for λ s that give largest probability to held-out set

Backoff

- Suppose you want:

$P(\text{pancakes} | \text{delicious soufflé})$

- If the trigram probability is 0, use the bigram

$P(\text{pancakes} | \text{soufflé})$

- If the bigram probability is 0, use the unigram

$P(\text{pancakes})$

- Complication: need to discount the higher-order n-gram so probabilities don't sum higher than 1 (e.g., Katz backoff)

Stupid Backoff

Backoff without discounting (not a true probability)

$$S(w_i|h) = \begin{cases} \frac{\text{count}(h, w_i)}{\text{count}(h)}, & \text{if } \text{count}(h, w_i) > 0 \\ \lambda S(w_i|h^i), & \text{otherwise} \end{cases}$$

- h = current history (n-1 words)
- h' = shorter history
- λ = fixed backoff factor (≈ 0.4)