

N-gram Language Modeling

Introduction to N-gram Language Models

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

For example LM can predict the following sequence is more likely
all of a sudden I notice three guys standing on the
sidewalk

than does this same set of words in a different order:
on guys all I of notice sidewalk three a sudden standing
the

Why word prediction?

It's a helpful part of language tasks

- One reason is for generation: choosing contextually better words.

- 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 and not 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 Modeling (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:

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

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?

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

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

No! Too many possible sentences!

We'll never see enough data for estimating these

Notation

- To represent the probability of a particular random variable X_i taking on the value “the”, or $P(X_i = \text{“the”})$, we will use the simplification $P(\text{the})$.
- We’ll represent a sequence of n words either as $w_1 \dots w_n$ or $w_{1:n}$.
- Thus the expression $w_{1:n-1}$ means the string w_1, w_2, \dots, w_{n-1} ,
- but we’ll also be using the equivalent notation $w_{\leq n}$, which can be read as “all the elements of w from w_1 up to and including w_{n-1} ”.
- For the joint probability of each word in a sequence having a particular value $P(X_1 = w_1, X_2 = w_2, X_3 = w_3, \dots, X_n = w_n)$ we’ll use $P(w_1, w_2, \dots, w_n)$.

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

Reminder: The Chain Rule

Recall the definition of conditional probabilities

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

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

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

$P(\text{"The water of Walden Pond"}) =$

$P(\text{The}) \times P(\text{water}|\text{The}) \times P(\text{of}|\text{The water})$

$\times P(\text{Walden}|\text{The water of}) \times P(\text{Pond}|\text{The water of Walden})$

Markov Assumption

Simplifying assumption:



A. A. Markov (1886).

Andrei Markov

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

$$\approx P(\text{blue} | \text{beautifully})$$

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

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, we approximate each component in the product

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

Simplest case: Unigram model

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

Some automatically generated sentences from two different unigram models

To him swallowed confess hear both . Which . Of save on trail
for are ay device and rote life have

Hill he late speaks ; or ! a more to leg less first you enter

Months the my and issue of year foreign new exchange's September

were recession exchange new endorsed a acquire to six executives

Bigram model

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

Some automatically generated sentences from two different bigram models

Why dost stand forth thy canopy, forsooth; he is this palpable hit
the King Henry. Live king. Follow.

What means, sir. I confess she? then all sorts, he is trim, captain.

Last December through the way to preserve the Hudson corporation N.
B. E. C. Taylor would seem to complete the major central planners
one gram 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

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?

A nice clear paradigm that lets us introduce many of the important issues for **large language models**

- **training** and **test** sets
- the **perplexity** metric
- **sampling** to generate sentences
- ideas like **interpolation** and **backoff**

N-gram Language Modeling

Introduction to N-grams

N-gram Language Modeling

Estimating N-gram Probabilities

Estimating bigram probabilities

The Maximum Likelihood Estimate

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

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

An example

< s > I am Sam </ s >

< s > Sam I am </ s >

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

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

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

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

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 9332 sentences, $|V|=1446$, taken only 8 words

	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

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

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								
want								
to								
eat								
chinese								
food								
lunch								
spend								

Raw bigram probabilities

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

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

$P(I|< s >)$

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

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

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

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

$= .000031$

$P(i|< s >) = 0.25$

$P(\text{english}|\text{want})$
 $= 0.0011$

$P(\text{food}|\text{english})$
 $= 0.5$

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

What kinds of knowledge do N-grams represent?

- Some of the bigram probabilities encode some facts that we think of as strictly syntactic in nature,
- like the fact that what comes after eat is usually a noun or an adjective,
- or that what comes after to is usually a verb
- Others might be a fact about the personal assistant task, like the high probability of sentences beginning with the words I.
- And some might even be cultural rather than linguistic, like the higher probability that people are looking for Chinese versus English food.

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 ngrams

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)

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.
- Efficiency: quantize probabilities to 4-8 bits instead of 8-byte float

N-gram LM Toolkits

SRILM

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

KenLM

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

Evaluation and Perplexity

Language Modeling

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. 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

Intuition of perplexity as evaluation metric:
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 : Predicting upcoming words



Claude Shannon

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 : 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 : 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 5: the inverse

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

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

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

Probability range is $[0,1]$, perplexity range is $[1,\infty]$

Minimizing perplexity is the same as maximizing probability

language model is evaluating a test set with two sentences.

Sentence 1: "< s > The cat sat on the mat < /s >"

Sentence 2: "< s > A dog chased the ball< /s >"

The model's probabilities for the words in each sentence are as follows:

Sentence 1: $P(\text{"T he"} | \text{< s >}) = 0.5$, $P(\text{"cat"} | \text{"T he"}) = 0.8$, $P(\text{"sat"} | \text{"T he cat"}) = 0.7$, $P(\text{"on"} | \text{"T he cat sat"}) = 0.9$, $P(\text{"the"} | \text{"T he cat sat on"}) = 0.6$, $P(\text{"mat"} | \text{"T he cat sat on the"}) = 0.7$

Sentence 2: $P(\text{"A"} | \text{< s >}) = 0.4$, $P(\text{"dog"} | \text{"A"}) = 0.6$, $P(\text{"chased"} | \text{"A dog"}) = 0.8$, $P(\text{"the"} | \text{"A dog chased"}) = 0.7$, $P(\text{"ball"} | \text{"A dog chased the"}) = 0.9$

Calculate the perplexity of the model on this test set.

$$P(S1)=P(\text{"The cat sat on the mat"}) \\ =0.5 \times 0.8 \times 0.7 \times 0.9 \times 0.6 \times 0.7 = 0.10584$$

$$P(S2)=P(\text{"A dog chased the ball"}) \\ =0.4 \times 0.6 \times 0.8 \times 0.7 \times 0.9 = 0.12096$$

$$P(\text{Test Set})=P(S1) \times P(S2) = 0.10584 \times 0.12096 = 0.012803$$

$$\text{Total words (N)} = 6+5=11$$

$$PP=(\text{Total Probability})^{-1/N}$$

$$PP=(0.012803)^{-1/11}$$

$$PP1.01286$$

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 $T = \text{"red red red red blue"}$

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

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:

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

$$= (.8 * .8 * .8 * .8 * .1)^{-1/5} = .04096^{-1/5} = .527^{-1} = 1.89$$

Holding test set constant:

Lower perplexity = better language model

Training 38 million words, test 1.5 million words, WSJ

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

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 training and 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.

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"

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

- 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

Sampling and Generalization

Language Modeling

The Shannon (1948) Visualization Method

Sample words from an LM



Claude Shannon

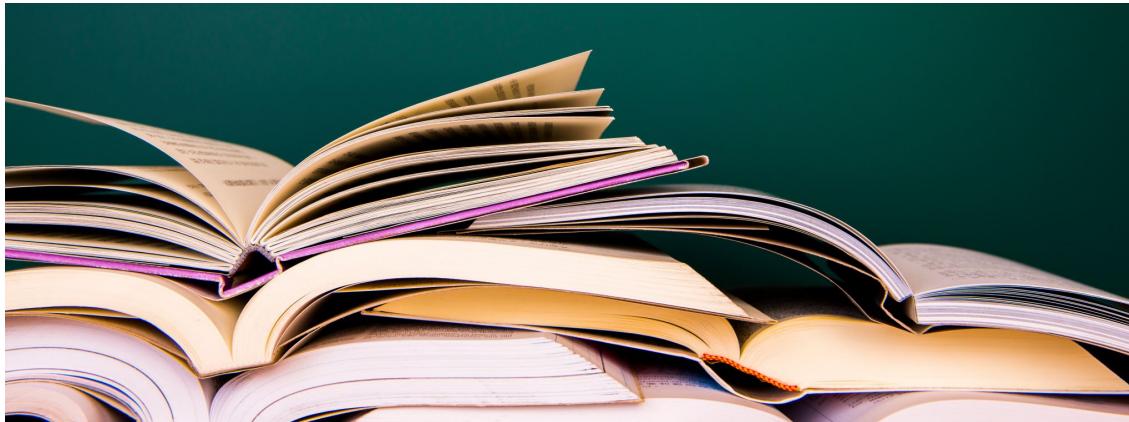
Unigram:

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

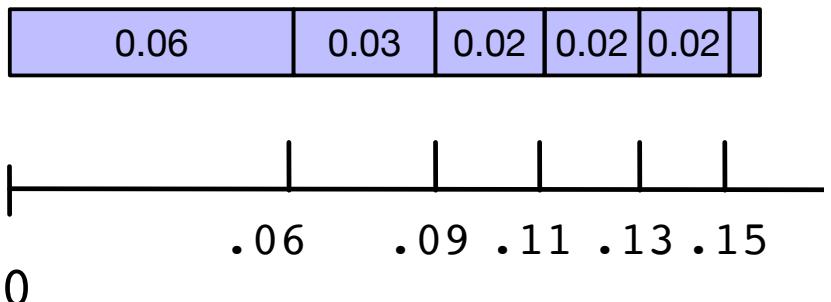


"Open a book at random and select a letter at random on the page. This letter is recorded. The book is then opened to another page and one reads until this letter is encountered. The succeeding letter is then recorded. Turning to another page this second letter is searched for and the succeeding letter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

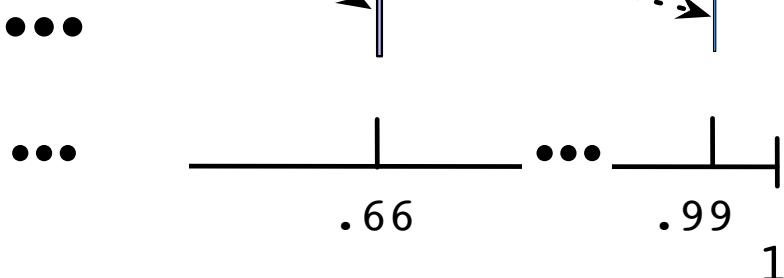
Sampling a word from a distribution



the of a to in



however ($p=.0003$) polyphonic
 $p=.0000018$



Visualizing Bigrams the Shannon Way

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

$\langle s \rangle$ I

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

I want

Now choose a random bigram (w, x)

want to

according to its probability $p(x|w)$

to eat

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

eat Chinese

Then string the words together

Chinese food

food $\langle /s \rangle$

I want to eat Chinese food

Approximating Shakespeare

1
gram

- To him swallowed confess hear both. Which. Of save on trail for are ay device and
rote life have
- Hill he late speaks; or! a more to leg less first you enter

2
gram

- Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live
king. Follow.
- What means, sir. I confess she? then all sorts, he is trim, captain.

3
gram

- Fly, and will rid me these news of price. Therefore the sadness of parting, as they say,
'tis done.
- This shall forbid it should be branded, if renown made it empty.

4
gram

- King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A
great banquet serv'd in;
- It cannot be but so.

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.

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

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

Zeros

Training set:

... ate lunch

... ate dinner

... ate a

... ate the

- Test set

- ... ate lunch

- ... ate breakfast

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

Zero probability 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)!

N-gram Language Modeling

Smoothing, Interpolation, and Backoff

The intuition of smoothing (from Dan Klein)

When we have sparse statistics:

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

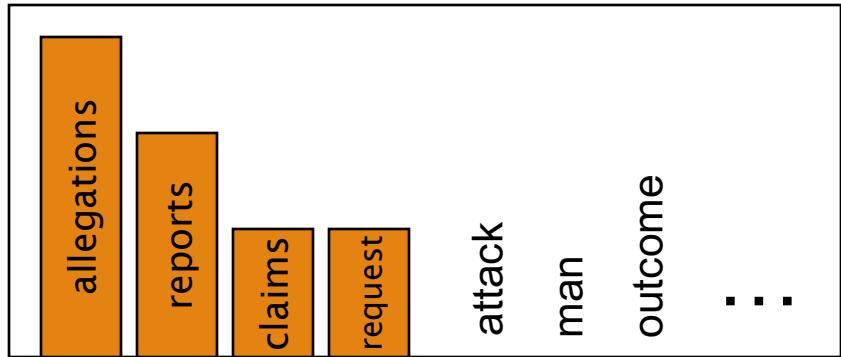
3 allegations

2 reports

1 claims

1 request

7 total



Steal probability mass to generalize better

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

2.5 allegations

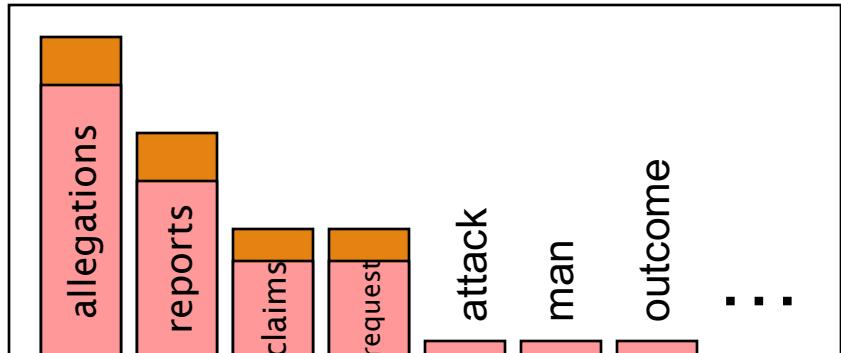
1.5 reports

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:

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

Add-1 estimate:

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

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

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

Reconstituted 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

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

$$\sum_i \lambda_i = 1$$

How to set λ s for interpolation?

Use a **held-out** corpus

Training Data

Held-Out
Data

Test
Data

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:

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

If the trigram probability is 0, use the bigram

$$P(\text{pancakes} \mid \text{soufflé})$$

If the bigram probability is 0, use the unigram

$$P(\text{pancakes})$$

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

why discounting is needed in backoff for n-gram model

Core intuition

- In a pure n-gram model, all observed n-grams get probability proportional to their counts, and unseen n-grams get probability 0.
- Backoff says: if an n-gram is unseen (count 0), use a lower-order model (e.g., backoff from trigram to bigram, then to unigram).
- If the higher-order probabilities are not reduced first, then the lower-order probabilities added during backoff are “extra” mass, so the conditional distribution over next words no longer sums to 1.

What discounting does

- Discounting reduces each non-zero higher-order count c to a smaller effective count c^*
- This reduction “frees” some total probability mass, which is then assigned (through a backoff weight like $\alpha(h)$) to the lower-order distribution used for unseen or low-count events.

Example

Consider a tiny bigram (2-gram) model with vocabulary {I, eat, food}. The goal is to model $P(\text{next word}|\text{current word})$.

Assume the training corpus gives these bigram counts:

Count(I, eat)=10

Count(I, food)=0

Count(eat, food)=5

Unigram counts for histories:

Count(I)=10 (only followed by “eat” in data)

Count(eat)=5 (only followed by “food” in data)

Say unigram probabilities are:

$$P(\text{eat})=0.3, P(\text{I})=0.3, P(\text{food})=0.4$$

Now,

$$P(\text{eat}|\text{I})=10/10=1$$

$$P(\text{food}|\text{I})=0/10=0$$

$$P(\text{food}|\text{eat})=5/5=1$$

What naïve backoff would try to do

Suppose we decide:

If a bigram has count 0, back off to unigram $P(\text{food})$.

If a bigram has count > 0 , keep using bigram probability.

Then the conditional distribution after “I” would be:

For seen bigram: $P(\text{eat}|I)=1$ (from bigram)

For unseen bigram “I food”: add $P(\text{food})=0.4$ via backoff

Now the total probability after “I” is

$$1+0.4=1.4>1 \quad 1+0.4=1.4>1$$

Discounting to free probability mass

Take D=0.75. Then for history “I”:

Original count: $c(I, \text{eat})=10$

Discounted count: $c^*(I, \text{eat})=10-0.75=9.25$

Total original count for history “I” is 10, so the discounted bigram probability becomes

$P^*(\text{eat}|I)=9.25/10=0.925$, Now the sum over seen continuations after “I” is 0.925, not 1. The “missing” 0.075 is reserved for unseen bigrams after “I” (like “I food”)

Stupid Backoff

Backoff without discounting (not a true probability)

$$S(w_i | w_{i-k+1}^{i-1}) = \begin{cases} \frac{\text{count}(w_{i-k+1}^i)}{\text{count}(w_{i-k+1}^{i-1})} & \text{if } \text{count}(w_{i-k+1}^i) > 0 \\ 0.4S(w_i | w_{i-k+2}^{i-1}) & \text{otherwise} \end{cases}$$

$$S(w_i) = \frac{\text{count}(w_i)}{N}$$

N-gram Language Modeling

Interpolation and Backoff