

CS 2731

# Introduction to Natural Language Processing

Session 5: N-gram language models

---

Michael Miller Yoder

September 10, 2025



University of  
Pittsburgh

School of Computing and Information

# Quiz

- Go to **Quizzes > Quiz 09-10** on Canvas
- You have until **2:45pm** to complete it
- Allowed resources
  - Textbook
  - Your notes (on a computer or physical)
  - Course slides and website
- Resources not allowed
  - Generative AI
  - Internet searches

# Course logistics

- Homework 1 will be released today or tomorrow
- [Project idea form](#) is **due tomorrow, Thu Sep 11**
  - You will be able to submit any project ideas that you're interested in: from the [example list](#) or any you have on your own
  - It's fine to incorporate your own research, there just needs to be an NLP component
  - You can submit multiple project ideas
- You will later choose from an anonymized list of project ideas on Project Match Day, next Wed Sep 17

# Overview: N-gram language models part 1

- Language modeling
- N-gram language models
- Estimating n-gram probabilities
- Perplexity and evaluating language models
- Coding activity: build your own n-gram language model!

# Structure of this course

## MODULE 1

### Prerequisite skills for NLP

text normalization, linear alg., prob., machine learning

## MODULE 2

statistical machine learning

n-grams

language modeling  
text classification

## MODULE 3

## MODULE 4

language modeling  
text classification  
sequence labeling

## MODULE 5

### NLP applications and ethics

machine translation, chatbots, information retrieval, bias

# Introduction to language models

---

# Language Models Estimate the Probability of Sequences

Which of these sentences would you be more likely to observe in an English corpus?

- Hugged I big brother my.
- I hugged my large brother.
- I hugged my big brother.



# Language Models Estimate the Probability of Sequences

Which of following word would be most likely to come after “David hates visiting New...”

- York
- California
- giggled





These are actually instances of the same problem: the language modeling problem!

# Language Modeling is Tremendously Useful

LMS (language models) are at the center of NLP today and have many different applications

- **Machine Translation**

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

- **Spelling Correction**

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

- **Text Input Methods**

$P(\text{i cant believe how hot you are}) > P(\text{i cant believe how hot you art})$

- **Speech Recognition**

$P(\text{recognize speech}) > P(\text{wreck a nice beach})$

# The Goal of Language Modeling

Compute the probability of a sequence of words/tokens/characters:

$$P(\mathbf{w}) = P(w_1, w_2, w_3, w_5, \dots, w_n)$$

$$P(\text{I, hugged, my, big, brother})$$

This is related to next-word prediction:

$$P(w_t | w_1 w_2 \dots w_{t-1})$$

$$P(\text{York} | \text{David, hates, going, to, New})$$

Do you compute either of these? Then you're in luck:

**You are a language model!**

# N-gram language models

---

# The Chain Rule Helps Us Compute Joint Probabilities

The definition of conditional probability is

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

which can be rewritten as

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

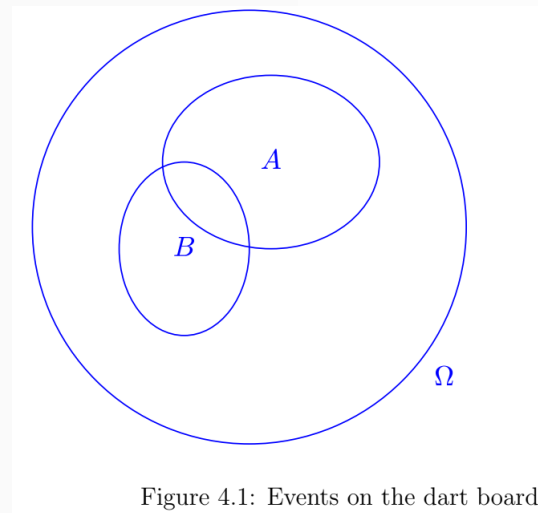


Figure 4.1: Events on the dart board

# The Chain Rule Helps Us Compute Joint Probabilities

If we add more variables, we see the following pattern:

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

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

which can be generalized as

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

**The Chain Rule!**

# The chain rule to compute the joint probability of words in a sentence

$$P(w_1, w_2, w_3, \dots, w_n) = \prod_i^n P(w_i | w_1 w_2 \dots w_{i-1})$$

$P(\text{now is the winter of our discontent}) =$   
 $P(\text{now}) \times P(\text{is}|\text{now}) \times$   
 $P(\text{the}|\text{now is}) \times P(\text{winter}|\text{now is the}) \times$   
 $P(\text{of}|\text{now is the winter}) \times$   
 $P(\text{our}|\text{now is the winter of}) \times$   
 $P(\text{discontent}|\text{now is the winter of our})$



# How Are We Estimating these Probabilities?

Could we just count and divide?

$$P(\text{discontent} | \text{now is the winter of our}) = \frac{\text{Count}(\text{now is the winter of our discontent})}{\text{Count}(\text{now is the winter of our})}$$

But this can't be a valid estimate! “now is the winter of our” is going to be very rare in corpora. It isn't going to be a good estimate of its true probability.



# This May not Seem Very Helpful

Is  $P(\text{discontent} | \text{now is the winter of our})$  really easier to compute than  $P(\text{now is the winter of our} | \text{discontent})$ ?

How can the chain rule help us? We can **cheat**.

# Markov Showed that You Could Make a Simplifying Assumption

One can approximate

$$P(\text{discontent}|\text{now is the winter of our})$$

by computing

$$P(\text{discontent}|\text{our})$$

or perhaps

$$P(\text{discontent}|\text{of our})$$



- We can obtain our estimate by only counting simpler things: “our discontent”, “discontent”, “of our”, etc
- N-gram language modeling is a generalization of this observation

# This assumption is the 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:

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

We will now walk through what this looks like for different values of  $k$ .

# The Unigram Model ( $k = 1$ )

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

The probability of a sequence is approximately the product of the probabilities of the individual words.

Some automatically generated sequences 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

What do you notice about them?

# The Bigram Model ( $k = 2$ )

If you condition on the previous word, you get the following:

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

Some examples generated by 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

Are these better?

# The Trigram Model

The trigram model is just like the bigram model, only with a larger  $k$ :

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx P(w_i | w_{i-2} w_{i-1})$$

The output of a trigram language model is generally **much** better than that of a bigram model **provided the training corpus is large enough**. Why do you need a larger corpus to train a trigram corpus than a bigram or unigram corpus?

# N-gram models have trouble with long-range dependencies

In general, n-gram models are very impoverished models of language. For example, language has relationships that span many words:

- The **students** who worked on the assignment for three hours straight **\*is/are** finally resting.
- The **teacher** who might have suddenly and abruptly met students **is/\*are** tall.
- Violins are easy to mistakenly think you can learn to play **\*them/quickly**.

# Ngram LMs Are Often Adequate

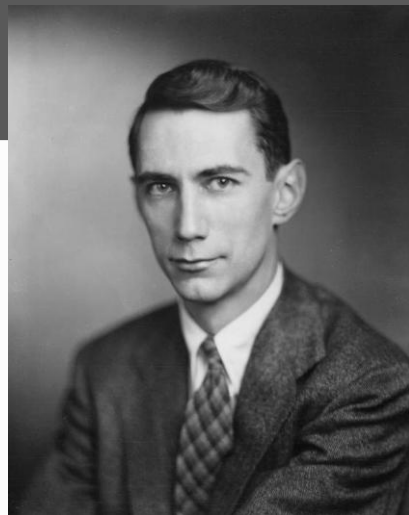
Nevertheless, for many applications, ngram models are good enough (and they're super fast and efficient)



# Sampling sentences from language models

---

# The Shannon Visualization Method



- Choose a random bigram ( $\langle s \rangle$ ,  $w$ ) according to its probability
- Now choose a random bigram ( $w$ ,  $x$ ) according to its probability
- 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

# Estimating n-gram probabilities

---

# Estimating bigram probabilities with the maximum likelihood estimate (MLE)

MLE for bigram probabilities can be computed as:

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

which we will sometimes represent as

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

# An example

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

<s> I am Sam </s>

<s> Sam I am </s>

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

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

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

$$P(\text{am} | \text{I}) =$$

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

$$P(\text{Sam} | \text{am}) =$$

$$P(\text{do} | \text{I}) =$$

## 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

# 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

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

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

$$\begin{aligned} & \times P(\text{want} | \text{I}) \\ & \times P(\text{english} | \text{want}) \\ & \times P(\text{food} | \text{english}) \\ & \times P(</s> | \text{food}) \end{aligned}$$

$$= .000031$$

# Multiplication Considered Harmful

Doing computation in log space is preferred for language models

- **Avoid underflow** Multiplying small probabilities by small probabilities results in *very small* numbers, which is problematic
- **Optimize computation** Addition is cheaper than multiplication

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

# Perplexity and evaluating language models

---

# The Evaluation Process for ML Models

The goal of LM evaluation:

- Does our model prefer good sentences to bad sentences?
- Specifically, does it assign higher probabilities to the good/grammatical/frequently observed ones and lower probabilities to the bad/ungrammatical/seldom observed ones?

In ML evaluation, we divide our data into three sets: **train**, **dev**, and **test**.

- We train the model's parameters on the **train** set
- We tune the model's hyperparameters (if appropriate) on the **dev** set (which should not overlap with the **train** set)
- We test the model on the **test** set, which should not overlap with **train** or **dev**

An **evaluation metric** tells us how well our model has done on **test**.

# We Can Evaluate Models Intrinsically or Extrinsically

- **Extrinsic Evaluation** means asking how much the model contributes to a larger task or goal. We may evaluate an LM based on how much it improves machine translation over a BASELINE.
- **Intrinsic Evaluation** means measuring some property of the model directly. We may quantify the probability that an LM assigns to a corpus of text.

In general, EXTRINSIC EVALUATION is better, but more expensive and time-consuming.

# Extrinsic Evaluation of LMs

## Best evaluation for comparing models A and B

- Put each model in a task (spelling corrector, speech recognizer, MT system)
- Run the task, get an accuracy for A and for B
  - How many misspelled words corrected properly?
  - How many sentences translated correctly?
- Compare scores for A and B

**This takes a lot of time to set up and can be expensive to carry out.**

# Perplexity is an intrinsic metric for language modeling

Perplexity evaluates the probability assigned by a model **to a collection of test documents, controlling for length** and is, thus, useful for evaluating LMs.

A better model of a text is one which assigns a higher probability to words that actually occur in the test set. This will result in **lower** perplexity.



However:

- It is a rather crude instrument
- It sometimes correlates only weakly with performance on downstream tasks
- It's only useful for pilot experiments
- But it's cheap and easy to compute, so it's important to understand

# Deriving Perplexity for Bigrams

$$PP(w) = P(w_1 w_2 \dots w_n)^{-\frac{1}{n}}$$

Definition

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

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

Chain Rule

$$= \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i)}}$$

For Unigrams

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

For Bigrams

**To minimize perplexity is to maximize probability!**



In general, a lower perplexity implies a better model.

# Lower perplexity = better model

Training 38 million words, test 1.5 million words,  
WSJ

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

# The problem of zeros

---

# The Perils of Overfitting

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

- In real life, it often doesn't
- 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

# N-grams in the test set that weren't in the training set

Suppose our bigram LM, trained on Twitter, reads a document by the philosopher Wittgenstein:

*Whereof one cannot speak, thereof one must be silent.*

This contains the bigrams: whereof one, one cannot, cannot speak, speak [comma], [comma] thereof, thereof one, one must, must be, be silent.

Suppose “whereof one” never occurs in the training corpus (**train**) but whereof occurs 20 times. According to MLE, it's probability is

$$P(\text{one}|\text{whereof}) = \frac{c(\text{whereof, one})}{c(\text{whereof})} = \frac{0}{20} = 0$$

The probability of the sentence is the **product** of the probabilities of the bigrams. What happens if one of the probabilities is zero?

# Laplace and Lidstone smoothing

---

# The intuition of smoothing

When we have sparse statistics:

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

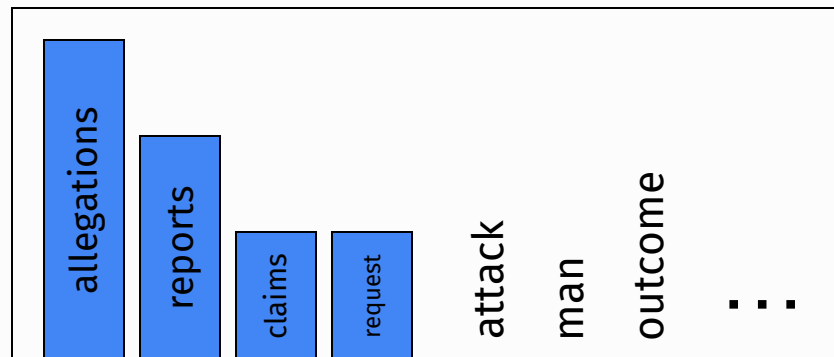
3 allegations

2 reports

1 claims

1 request

7 total



Steal probability mass to generalize better

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

2.5 allegations

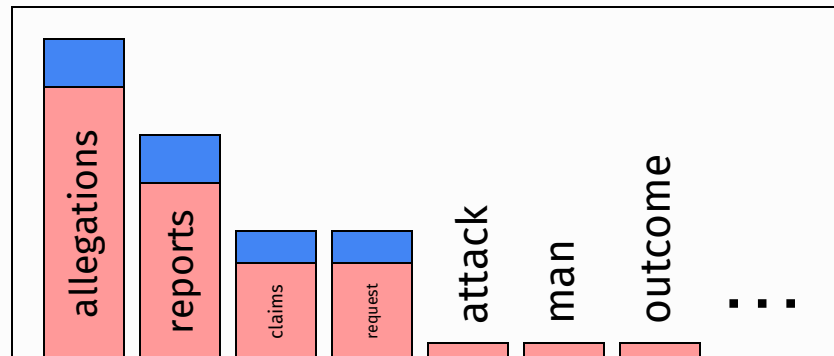
1.5 reports

0.5 claims

0.5 request

2 other

7 total



# Laplace smoothing: Pretending that we saw each word once more

$$\text{MLE estimate } P_{MLE}(w_i|w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$\text{Add-1 estimate } P_{Add-1}(w_i|w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + |V|}$$

Where  $V$  is the vocabulary of the corpus.



# Laplace smoothing is too blunt

Problem: A large dictionary makes rare words too probable.

Solution: instead of adding 1 to all counts, add  $k < 0$ .

How to choose  $k$ ?

# How to choose k?

## Add-0.001 Smoothing

Doesn't smooth much

xya	1	1/3	1.001	0.331
xyb	0	0/3	0.001	0.0003
xyc	0	0/3	0.001	0.0003
xyd	2	2/3	2.001	0.661
xye	0	0/3	0.001	0.0003
...				
xyz	0	0/3	0.001	0.0003
Total xy	3	3/3	3.026	1

72

# How to choose $k$ ?

- Hyperparameter!
  - Try many  $k$  values on dev data and choose  $k$  that gives the lowest perplexity
  - Report result on test data
- Could tune this at the same time as  $n$  in  $n$ -gram LM

# Coding activity: build your own n-gram LMs

---

# Clickbait n-grams and classification on JupyterHub

- [Click on this nbgitpuller link](#)
- Open `session5_ngram_lm.ipynb`

*Questions?*