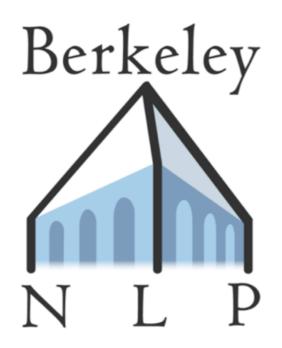
# N-gram Language Models



EECS 183/283a: Natural Language Processing

### N-Grams



- An *n*-gram is a sequence of *n* tokens
- ullet Given a vocabulary  ${\mathcal V}$ , the set of all possible n-grams is  ${\mathcal V}^n$

#### A language model assigns a

• • •

### Autoregressive Language Models



$$p(\overline{x}) = \prod_{i=1}^{|\overline{x}|} p(x_i \mid x_1, \dots, x_{i-1})$$

$$= p(x_1) p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$$

probability that probability that the first word is  $x_1$  the second word is  $x_2$ , given that the first word is  $x_1$ 

probability that the sentence ends after the sequence  $\langle x_1, \ldots, x_{n-1} \rangle$ 

#### Core modeling challenge:

How do we compute these conditional probabilities?

# Autoregressive N-Gram Language Models



$$p(\overline{x}) = \prod_{i=1}^{|\overline{x}|} p(x_i \mid x_1, \dots, x_{i-1})$$

$$= p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$$

$$p(X_i = x) = p(x \mid x_1, \dots, x_{i-1})$$

Challenge: how do we parameterize this if there can be arbitrarily many previous tokens?

#### Let's make a Markov assumption:

The probability of word at index i only depends on the n - 1 words that came before it

# Autoregressive N-Gram Language Models



$$p(\overline{x}) = \prod_{i=1}^{|\overline{x}|} p(x_i \mid x_1, \dots, x_{i-1})$$

$$= p(x_1)p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1})$$

$$p(X_i = x) = p(x \mid x_1, \dots, x_{i-1})$$

$$\approx p(x \mid x_{i-n+1}, \dots, x_{i-1})$$

#### Let's make a Markov assumption:

The probability of word at index i only depends on the n - 1 words that came before it

# Autoregressive N-Gram Language Models



$$\begin{split} p(\overline{x}) &= \prod_{i=1}^{|\overline{x}|} p(x_i \mid x_1, \dots, x_{i-1}) \\ &= p(x_1) p(x_2 \mid x_1) \dots p(x_n \mid x_1, \dots, x_{n-1}) \\ p(X_i = x) &= p(x \mid x_1, \dots, x_{i-1}) \\ &\approx p(x \mid \underbrace{x_{i-n+1}, \dots, x_{i-1}}) \\ &\text{Preceding (n-1)-gram} \end{split}$$

$$pprox rac{C(x_{i-n+1},\ldots,x_{i-1},x)}{C(x_{i-n+1},\ldots,x_{i-1})}$$
 count of n-gram

# N-Gram Language Models



- Now, all we need to model is n-gram and (n-1)-gram probabilities!
- We can do this by counting n-gram occurrences in the wild
- Simplest n-gram language model: n = 1 (unigrams, aka bag of words)

total number of words in corpus 
$$p(X_i = x) \approx p(x) \in \Delta^{\mathcal{V}} \text{ across corpus}$$

$$= \frac{C(x)}{\sum_{\overline{x} \in \mathcal{D}} |\overline{x}|}$$

$$\text{words in corpus}$$

### Unigram Language Model



$$p(\overline{x}) = \prod_{i=1}^{|\overline{x}|} p(x_i \mid x_1, \dots, x_{i-1})$$

$$\approx \prod_{i=1}^{|\overline{x}|} p(x_i)$$

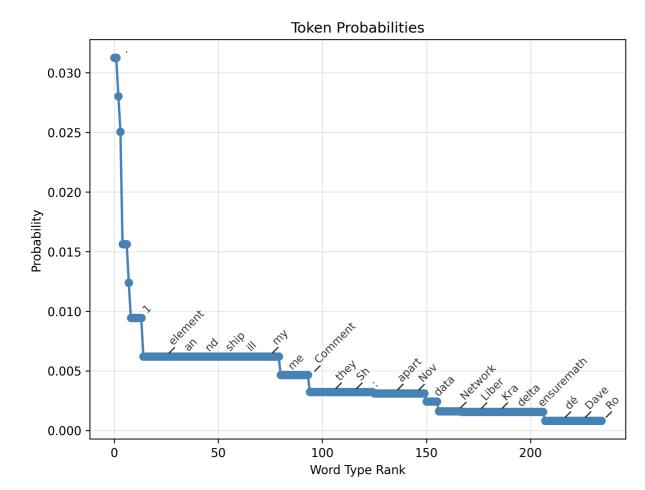
p(yellow suitcase and red hat) = p(red suitcase and yellow hat)

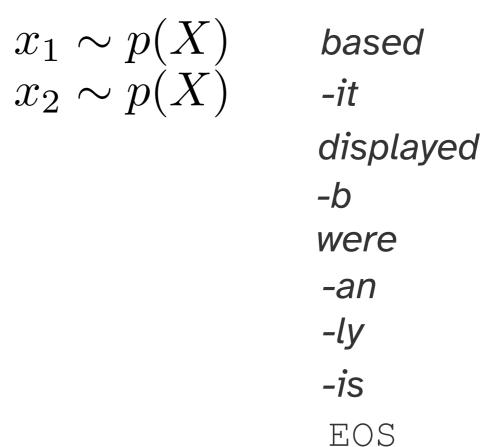
Word order does not matter!
This is why it's called "bag of words"

### Unigram Language Model



p(X) (unigram probabilities)





basedit displayedb wereanlyis

What are the parameters of this model?

How many parameters are there?



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

$$= \frac{C(x_{i-1}, x_i)}{C(x_{i-1})}$$

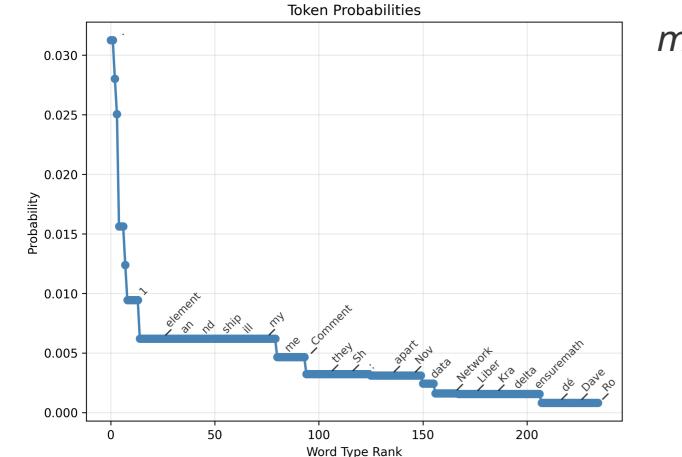
$$p(\overline{x}) \approx \prod_{i=1}^{|\overline{x}|} \frac{C(x_{i-1}, x_i)}{C(x_{i-1})}$$



 First sample a start token using the empirical first-token distribution:

$$p(X_1 = x) = \frac{C(\mathrm{BOS}, x)}{C(\mathrm{BOS})} = \frac{C(\mathrm{BOS}, x)}{|\mathcal{D}|}$$

$$p(X \mid \mathtt{BOS})$$



my

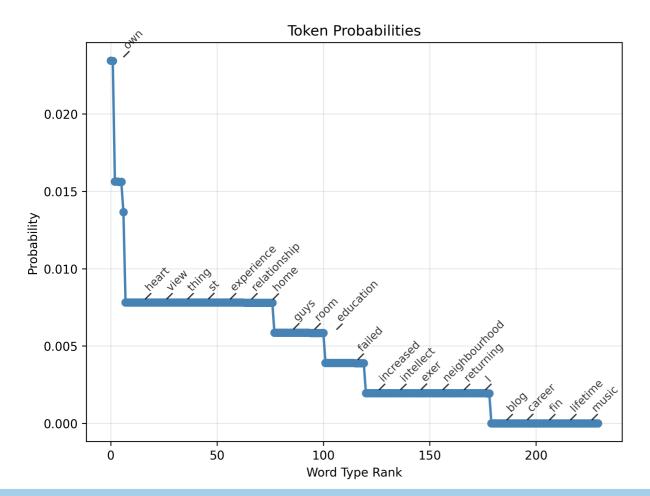
Infini-gram (Liu et al. 2024) stats, vibecoded with Claude



Then sample conditioned on the previous word:

$$p(X_2 = x) = \frac{C(\text{my}, x)}{C(\text{my})}$$

$$p(X \mid my)$$



my mother

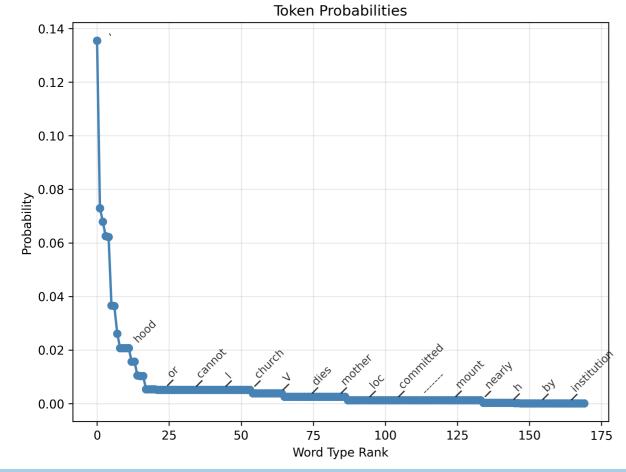
Infini-gram (Liu et al. 2024) stats, vibecoded with Claude



Then sample conditioned on the previous word:

$$p(X_3 = x) = \frac{C(\text{mother}, x)}{C(\text{mother})}$$

$$p(X \mid \text{mother})$$



my mother and

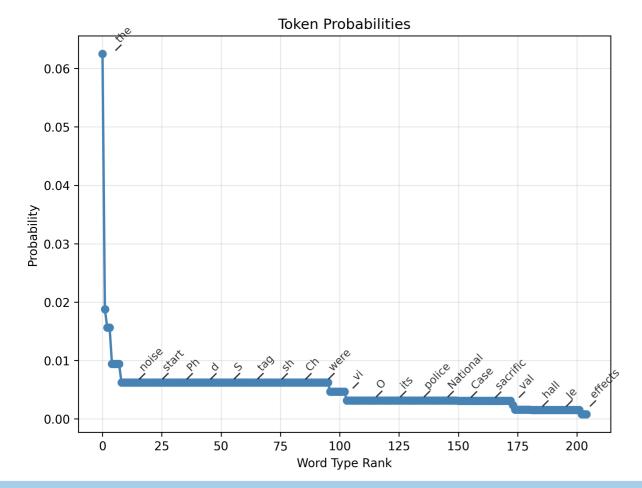
Infini-gram (Liu et al. 2024) stats, vibecoded with Claude



Then sample conditioned on the previous word:

$$p(X_4 = x) = \frac{C(\text{and}, x)}{C(\text{and})}$$

$$p(X \mid \text{and})$$



my mother and my

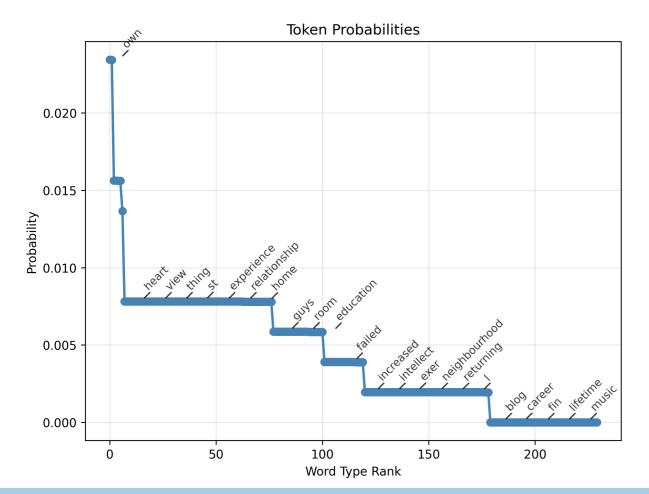
Infini-gram (Liu et al. 2024) stats, vibecoded with Claude



Then sample conditioned on the previous word:

$$p(X_5 = x) = \frac{C(\text{my}, x)}{C(\text{my})}$$

$$p(X \mid my)$$



my mother and my mother

•••

Infini-gram (Liu et al. 2024) stats, vibecoded with Claude



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

n=1		
never		
Comment		
has		
in		
•		
44		
t		
view		
С		
never		
Comment		
has in . "t		
view C		



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

n=1	n=2		
never	view		
Comment	-find		
has	а		
in	place		
•	of		
"	а		
t	human		
view	intelligence		
С	brief		
never	viewfind a		
Comment	place of a		
has in . "t	human		
view C	intelligence		
VIEW C	brief		



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

n=1	n=2	n=3	
never	view	were	
Comment	-find	е	
has	а	-colog	
in	place	-ical	
•	of	topics	
46	а	related	
t	human	Т	
view	intelligence	-weet	
С	brief	niet	
never	viewfind a	were	
Comment	place of a	ecological	
has in . "t	human	topics related	
view C	intelligence brief	to Tweet niet	



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

n=1	n=2	n=3	n=4	
never	view	were	dropped	
Comment	-find	е	her	
has	а	-colog	off	
in	place	-ical	at	
•	of	topics	her	
46	а	related	achiev	
t	human	Т	-ement	
view	intelligence	-weet	Canadian	
С	brief	niet	Resource	
never	viewfind a	were	dropped her off at	
Comment	place of a	ecological	her achievement	
has in . "t	human	topics related	Canadian	
view C	intelligence brief	to Tweet niet	1	



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

n=1	n=2	n=3	n=4	n=10
never	view	were	dropped	Liber
Comment	-find	е	her	-als
has	а	-colog	off	9
in	place	-ical	at	Third
•	of	topics	her	Part
46	a	related	achiev	-ies
t	human	Т	-ement	,
view	intelligence	-weet	Canadian	Left
С	brief	niet	Resource	_
never	viewfind a	were	dropped her off at	
Comment	place of a	ecological	her achievement	Liberals, Third
has in . "t	human	topics related	Canadian	Parties, Left-
view C	intelligence brief	to Tweet niet	Resource	

Sampled from Infini-gram (Liu et al. 2024)



As *n* increases, we get more fluent text

But also more sparsity:  $p(X_{10} = -) = \frac{C(\text{Liberals, Third Parties, Left-})}{C(\text{Liberals, Third Parties, Left})} = \frac{417}{418}$  (in a corpus of 1.4T tokens!)

n=1	n=2	n=3	n=4	n=10
never	view	were	dropped	Liber
Comment	-find	е	her	-als
has	а	-colog	off	,
in	place	-ical	at	Third
•	of	topics	her	Part
"	а	related	achiev	-ies
t	human	Т	-ement	,
view	intelligence	-weet	Canadian	Left
С	brief	niet	Resource	_
never Comment has in . "t view C	viewfind a place of a human intelligence brief	were ecological topics related to Tweet niet		Liberals, Third Parties, Left-

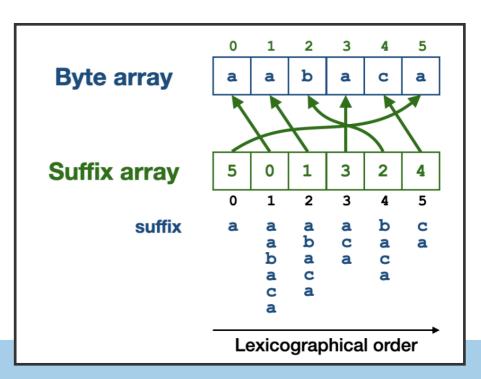
Sampled from Infini-gram (Liu et al. 2024)

### Storage Size



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

- For an *n*-gram language model, we need to store counts for:
  - All sequences of length n (  $\mathcal{V}^n$ )
  - All sequences of length n 1 (  $\mathcal{V}^{n-1}$  )
- There are actually more efficient ways to "store" *n*-gram models! See infini-gram (Liu et al. 2024)





$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

#### Missing data (sparsity)

- What if the count of the target n-gram is 0?
  - Solution: add a small number to the count for every ngram (aka "smoothing")



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

#### Missing data (sparsity)

- What if the count of the target n-gram is 0?
  - Solution: add a small number to the count for every ngram (aka "smoothing")
- What if our n-1-gram prefix has a count of 0?
  - Solution: condition on a shorter *n*-gram prefix (e.g., the previous *n*-2, or *n*-3, etc.) instead (aka "backoff")



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

#### **Storage space**

- What if the count of the target n-gram is 0?→ smoothing
- What if our n-1-gram prefix has a count of  $0? \rightarrow backoff$
- We need to store the counts for all n-grams we've seen in our corpus at worst, exponential wrt  $|\mathcal{V}|$



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

#### No notion of similarity

- What if the count of the target n-gram is 0?→ smoothing
- What if our n-1-gram prefix has a count of  $0? \rightarrow backoff$
- ullet Storage is at worst exponential wrt  $|\mathcal{V}|$
- Can't learn anything from the counts of n-grams containing similar words

$$p(\text{bike} \mid \text{I bought a}) \approx p(\text{bicycle} \mid \text{I purchased a})$$



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

#### Cannot condition on context with intervening words

- What if the count of the target n-gram is 0?→ smoothing
- What if our n-1-gram prefix has a count of  $0? \rightarrow backoff$
- ullet Storage is at worst exponential wrt  $|\mathcal{V}|$
- No notion of similarity

$$p(\text{Smith} \mid \text{Dr. Jane}) \approx p(\text{Smith} \mid \text{Dr. John}) \approx p(\text{Smith} \mid \text{Dr. } ---)$$



$$p(X_i = x) \approx \frac{C(x_{i-n+1}, \dots, x_{i-1}, x)}{C(x_{i-n+1}, \dots, x_{i-1})}$$

#### Without a big n, cannot handle long-distance dependencies

- What if the count of the target n-gram is 0?→ smoothing
- What if our n-1-gram prefix has a count of  $0? \rightarrow backoff$
- ullet Storage is at worst exponential wrt  $|\mathcal{V}|$
- No notion of similarity
- Intervening words

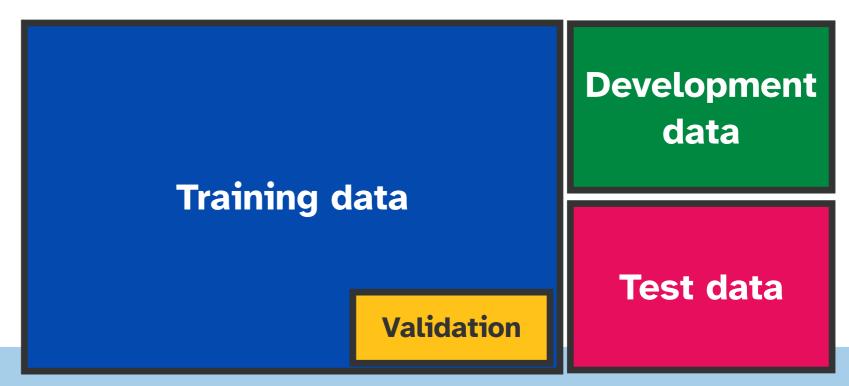
my cat Pepper, who is black with green eyes, likes to run on the cat wheel.



# Evaluating a Language Model



- Let's say we have a language model that can give us a probability of any text  $p_{\theta}(\overline{x})$
- We created this language model using a corpus  $\mathcal{D} = \{\overline{x}_i\}_{i=1}^M$
- We care how well this generalizes to some held-out dataset  $\mathcal{D}'$



### Measures of Fit



• Likelihood: probability of the data under our model

$$\prod_{i=1}^{M} p_{\theta}(\overline{x}_i)$$

### Measures of Fit



Likelihood: probability of the data under our model

$$\prod_{i=1}^{M} p_{\theta}(\overline{x}_i)$$

• Negative log likelihood (fixes float underflow) M

$$-\sum_{i=1}^{M} \log p_{\theta}(\overline{x}_i) = -\sum_{i=1}^{M} \sum_{j=1}^{N} \log p_{\theta}(x_j^i \mid x_1^i, \dots, x_{j-1}^i)$$

### Measures of Fit



Likelihood: probability of the data under our model

$$\prod_{i=1}^{M} p_{\theta}(\overline{x}_i)$$

• Negative log likelihood (fixes float underflow) M

$$-\sum_{i=1}^{M} \log p_{\theta}(\overline{x}_i) = -\sum_{i=1}^{M} \sum_{j=1}^{N} \log p_{\theta}(x_j^i \mid x_1^i, \dots, x_{j-1}^i)$$

 Perplexity: inverse probability of data, normalized by number of tokens in the dataset

$$= \exp\left(-\frac{1}{\sum_{i=1}^{M} |\overline{x}^i|} \sum_{i=1}^{M} \sum_{j=1}^{N} \log p_{\theta}(x_j^i \mid x_1^i, \dots, x_{j-1}^i)\right)$$