

N-grams & Smoothing

Chapter 4

Overview

- Recap of N-gram Modeling
 - Markov Assumption
 - How to Create Relative Frequency Tables?
 - How to Compute Probability of an Utterance
- How to Generate Sentences
 - Methods, Examples, Summarized Steps
- Zeros & Sparsity
 - Zeros
 - Sparsity
 - Recap MLE
- Smoothing
 - Zipfian Distribution
 - Laplace Smoothing (Add-1 smoothing)
 - Good-Turing
 - Backoff
 - Smoothing for Web-scale N-grams
 - Smoothing Summary

Markov assumption

- Estimate the conditional probability of the next word without looking too far in the past

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

Recap of N-gram Modeling

Bigram:
$$P(x_1^n) = \prod_{k=1}^n P(x_k | x_1^{k-1}) = \prod_{k=1}^n P(x_k | x_{k-1})$$

Trigram:
$$P(x_1^n) = \prod_{k=1}^n P(x_k | x_1^{k-1}) = \prod_{k=1}^n P(x_k | x_{k-2}^{k-1})$$

4-gram:
$$P(x_1^n) = \prod_{k=1}^n P(x_k | x_1^{k-1}) = \prod_{k=1}^n P(x_k | x_{k-3}^{k-1})$$

Recap of Relative Frequency

For **N**-gram, $P(x_1 \dots x_n) = \prod_{k=1}^n P(x_k | x_1 \dots x_{k-1})$
 $\approx \prod_{k=1}^n P(x_k | x_{k-N+1}, \dots, x_{k-1})$

$$P(x_1^n) = \prod_{k=1}^n P(x_k | x_1^{k-1}) = \prod_{k=1}^n P(x_k | x_{k-N+1}^{k-1})$$

$$P(x_k | x_{k-N+1}^{k-1}) = \frac{\text{freq}(x_{k-N+1}^{k-1} x_k)}{\text{freq}(x_{k-N+1}^{k-1})}$$

this ratio is called relative frequency

Recap of N-gram Modeling Example

Let's see the first simple example from (Dr. Seuss' story):

<s> I am Sam </s>

<s> Sam I am </s>

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

Now we are going to use the **Bigram** model

$$P(x_k | x_{k-1}) = \frac{\text{freq}(x_{k-1}x_k)}{\text{freq}(x_{k-1})}$$

$$P(\text{ am} | \text{ I}) = P(\text{ I am}) / P(\text{ I}) = 2/3$$

$$P(\text{ I} | \text{ <s>}) = 2/3$$

$$P(\text{ Sam} | \text{ <s>}) = 1/3$$

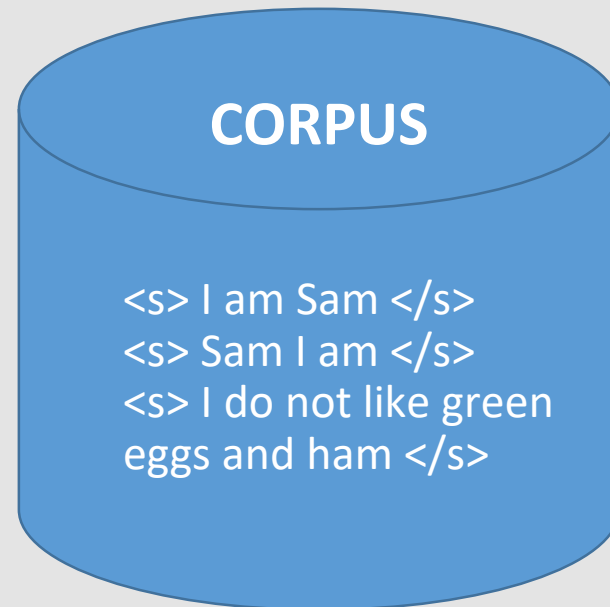
$$P(\text{ </s>} | \text{ Sam}) = 1/2$$

$$P(\text{ Sam} | \text{ am}) = 1/2$$

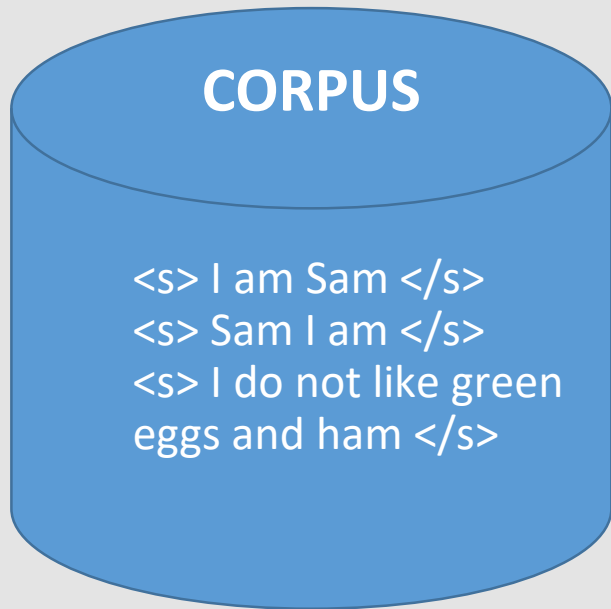
$$P(\text{ do} | \text{ I}) = 1/3$$

How to Create Relative Frequency Tables?

Example:



How to Create Relative Frequency Tables?



Unigram Table

Bigram Table

How to Create Relative Frequency Tables?

Unigram Table

I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
3	2	2	1	1	1	1	1	1	1	3	3

Bigram Table

	I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
I		2		1								
am			1									1
Sam	1											1
do					1							
not						1						
like							1					
green								1				
eggs									1			
and										1		
ham												1
<s>	2		1									
</s>												

CORPUS

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

How to Create Relative Frequency Tables?

Unigram Table

I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
3	2	2	1	1	1	1	1	1	1	3	3

Bigram Table

	I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
I	0	2	0	1	0	0	0	0	0	0	0	0
am	0	0	1	0	0	0	0	0	0	0	0	1
Sam	1	0	0	0	0	0	0	0	0	0	0	1
do	0	0	0	0	1	0	0	0	0	0	0	0
not	0	0	0	0	0	1	0	0	0	0	0	0
like	0	0	0	0	0	0	1	0	0	0	0	0
green	0	0	0	0	0	0	0	1	0	0	0	0
eggs	0	0	0	0	0	0	0	0	1	0	0	0
and	0	0	0	0	0	0	0	0	0	1	0	0
ham	0	0	0	0	0	0	0	0	0	0	0	1
<s>	2	0	1	0	0	0	0	0	0	0	0	0
</s>	0	0	0	0	0	0	0	0	0	0	0	0

CORPUS

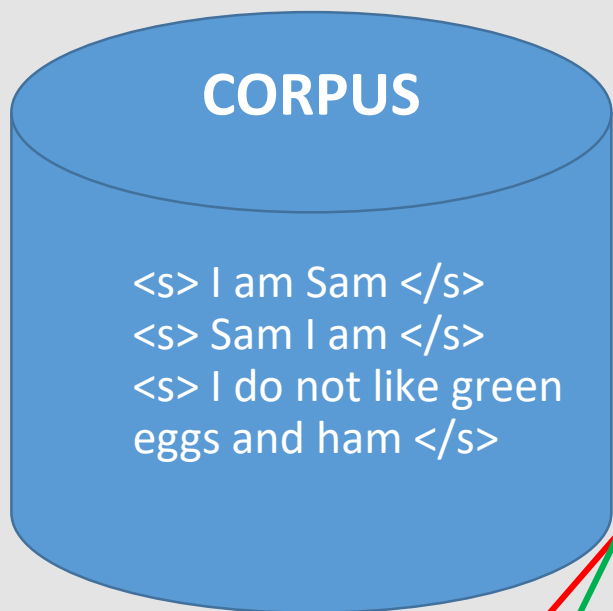
<s> I am Sam </s>
 <s> Sam I am </s>
 <s> I do not like green
 eggs and ham </s>

$$\text{Freq}(\text{am} | \text{I}) = \text{Freq}(\text{I am}) / \text{Freq}(\text{I}) =$$



Relative Frequencies →

How to Create Relative Frequency Tables?



$$\text{Freq}(\text{am} | \text{I}) = \text{Freq}(\text{I am}) / \text{Freq}(\text{I}) = 2/3$$

Unigram Table

Bigram Table

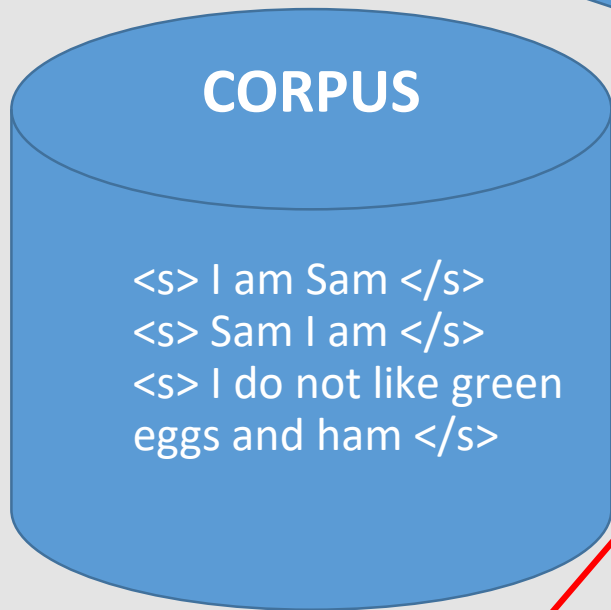
	I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
	3	2	2	1	1	1	1	1	1	1	3	3

	I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
I	0	2	0	1	0	0	0	0	0	0	0	0
am	0	0	1	0	0	0	0	0	0	0	0	1
Sam	1	0	0	0	0	0	0	0	0	0	0	1
do	0	0	0	0	1	0	0	0	0	0	0	0
not	0	0	0	0	0	1	0	0	0	0	0	0
like	0	0	0	0	0	0	1	0	0	0	0	0
green	0	0	0	0	0	0	0	1	0	0	0	0
eggs	0	0	0	0	0	0	0	0	1	0	0	0
and	0	0	0	0	0	0	0	0	0	1	0	0
ham	0	0	0	0	0	0	0	0	0	0	0	1
<s>	2	0	1	0	0	0	0	0	0	0	0	0
</s>	0	0	0	0	0	0	0	0	0	0	0	0

How to Create Relative Frequency Tables?

Unigram Table
Relative Frequencies

I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
3	2	2	1	1	1	1	1	1	1	3	3



$$\text{Freq}(\text{am} | \text{I}) = \text{Freq}(\text{I am}) / \text{Freq}(\text{I}) = 2/3$$

	I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
I	0	2/3	0	1/3	0	0	0	0	0	0	0	0
am	0	0	1/2	0	0	0	0	0	0	0	0	1/2
Sam	1/2	0	0	0	0	0	0	0	0	0	0	1/2
do	0	0	0	0	1	0	0	0	0	0	0	0
not	0	0	0	0	0	1	0	0	0	0	0	0
like	0	0	0	0	0	0	1	0	0	0	0	0
green	0	0	0	0	0	0	0	1	0	0	0	0
eggs	0	0	0	0	0	0	0	0	1	0	0	0
and	0	0	0	0	0	0	0	0	0	1	0	0
ham	0	0	0	0	0	0	0	0	0	0	0	1
<s>	2/3	0	1/3	0	0	0	0	0	0	0	0	0
</s>	0	0	0	0	0	0	0	0	0	0	0	0

How to Compute Probability of an Utterance?

- Given an n-gram language model

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

$$\begin{aligned} &P(<s>) * \\ &P(\text{I} | <s>) * \\ &P(\text{am} | \text{I}) * \\ &P(\text{Sam} | \text{am}) * \\ &P(</s> | \text{Sam}) \end{aligned}$$

This is called the *chain rule of probability*
using *the markov assumption*

$$P(<s> \mid \text{I am Sam } </s>) =$$

- $P(<s>)$ *

- $P(\text{I} \mid <s>)$ *

- $P(\text{am} \mid \text{I})$ *

- $P(\text{Sam} \mid \text{am})$ *

- $P(</s> \mid \text{Sam})$



Bigram Relative Frequency Table (from training set)

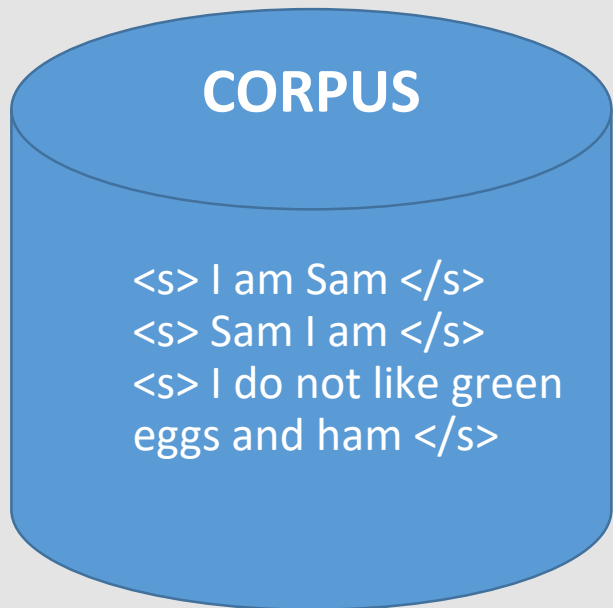
	I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
I	0	2/3	0	1/3	0	0	0	0	0	0	0	0
am	0	0	1/2	0	0	0	0	0	0	0	0	1/2
Sam	1/2	0	0	0	0	0	0	0	0	0	0	1/2
do	0	0	0	0	1	0	0	0	0	0	0	0
not	0	0	0	0	0	1	0	0	0	0	0	0
like	0	0	0	0	0	0	1	0	0	0	0	0
green	0	0	0	0	0	0	0	1	0	0	0	0
eggs	0	0	0	0	0	0	0	0	1	0	0	0
and	0	0	0	0	0	0	0	0	0	1	0	0
ham	0	0	0	0	0	0	0	0	0	0	0	1
<s>	2/3	0	1/3	0	0	0	0	0	0	0	0	0
</s>	0	0	0	0	0	0	0	0	0	0	0	0

How to Create Relative Frequency Tables?

$$P(\text{Unigram}) = \text{Freq}(\text{Unigram}) / N$$

N=20

I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
3/20	1/10	1/10	1/20	1/20	1/20	1/20	1/20	1/20	1/20	3/20	3/20



	I	am	Sam	do	not	like	green	eggs	and	ham	<s>	</s>
I	0	2/3	0	1/3	0	0	0	0	0	0	0	0
am	0	0	1/2	0	0	0	0	0	0	0	0	1/2
Sam	1/2	0	0	0	0	0	0	0	0	0	0	1/2
do	0	0	0	0	1	0	0	0	0	0	0	0
not	0	0	0	0	0	1	0	0	0	0	0	0
like	0	0	0	0	0	0	1	0	0	0	0	0
green	0	0	0	0	0	0	0	1	0	0	0	0
eggs	0	0	0	0	0	0	0	0	1	0	0	0
and	0	0	0	0	0	0	0	0	0	1	0	0
ham	0	0	0	0	0	0	0	0	0	0	0	1
<s>	2/3	0	1/3	0	0	0	0	0	0	0	0	0
</s>	0	0	0	0	0	0	0	0	0	0	0	0

Relative Frequencies →

How to Compute Probability of an Utterance?

- Given an n-gram language model

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

$$\begin{aligned} &P(<s>) * \\ &P(\text{I} | <s>) * \\ &P(\text{am} | \text{I}) * \\ &P(\text{Sam} | \text{am}) * \\ &P(</s> | \text{Sam}) \end{aligned}$$



This is called the *chain rule of probability*
using *the markov assumption*

$$= 3/20 * 2/3 * 2/3 * 1/2 * 1/2$$

Markov Assumption

- Estimate the conditional probability of the next word without looking too far in the past

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

For example, the sentence is “<s> I do not like green eggs and ham </s>”.

$P(\text{eggs} | \text{<s> I do not like green}) = P(\text{eggs} | \text{green})$ **Using a bigram model**

$P(\text{eggs} | \text{<s> I do not like green}) = P(\text{eggs} | \text{like green})$ **Using a trigram model**

$P(\text{eggs} | \text{<s> I do not like green}) = P(\text{eggs} | \text{not like green})$ **Using a 4-gram model**

etc ...

Relative Frequencies use **Two** Tables

N-gram table

and

(N-1)-gram table

How to Generate Sentences

Let's see BeRP example again...

i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
2533	927	2417	746	158	1093	341	278	3000	3000

Unigram table of raw frequencies

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	5	827	0	9	0	0	0	2	0	0
want	2	0	608	1	6	6	5	1	0	0
to	2	0	4	686	2	0	6	211	0	0
eat	0	0	2	0	16	2	42	0	0	34
chinese	1	0	0	0	0	82	1	0	0	23
food	15	0	15	0	1	4	0	0	0	12
lunch	2	0	0	0	0	1	0	0	0	9
spend	1	0	1	0	0	0	0	0	1	17
<start>	45	0	30	0	15	10	3	0	0	0
<end>	0	0	0	0	3	23	6	34	0	0

Bigram table of raw frequencies

P(I want to eat chinese food)?

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	0.002	0.33	0	0.0036	0	0	0	0.00079	0	0
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011	0	0
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087	0	0
eat	0	0	0.0027	0	0.021	0.0027	0.056	0	0	0.011
chinese	0.0063	0	0	0	0	0.52	0.0063	0	0	0.008
food	0.014	0	0.014	0	0.00092	0.0037	0	0	0	0.004
lunch	0.0059	0	0	0	0	0.0029	0	0	0	0.003
spend	0.0036	0	0.0036	0	0	0	0	0	1	0.006
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0
<end>	0	0	0	0	0.001	0.007	0.002	0.011	0	0

Relative Frequency Table

How to Generate Sentences

We know to **begin** the sentences
we want to use the **<start>** tag

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	0.002	0.33	0	0.0036	0	0	0	0.00079	0	0
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011	0	0
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087	0	0
eat	0	0	0.0027	0	0.021	0.0027	0.056	0	0	0.011
chinese	0.0063	0	0	0	0	0.52	0.0063	0	0	0.008
food	0.014	0	0.014	0	0.00092	0.0037	0	0	0	0.004
lunch	0.0059	0	0	0	0	0.0029	0	0	0	0.003
spend	0.0036	0	0.0036	0	0	0	0	0	1	0.006
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0
<end>	0	0	0	0	0.001	0.007	0.002	0.011	0	0

How to Generate Sentences

Sentences **begin** with the
<start> tag

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0

Now we are only interested in those words that
follow **<start>**
(the **non-zero** elements)

Why?

Because we are using our language model
(the **relative frequency table**)
to generate the words in our sentence

How to Choose?

If we pick the one with **the highest probability**, our sentences are **not** going to **change** very much

So

randomly pick one based on the **distribution**

	i	to	chinese	food	lunch
<start>	0.015	0.01	0.005	0.003	0.001

Which of the five choices do we choose

Randomly pick one based on its distribution

	i	to	chinese	food	lunch
<start>	0.015	0.01	0.005	0.003	0.001

↓ ↓ ↓ ↓ ↓

0.015 0.01 0.005 0.003 0.001

To do this we need to **normalize** **<start>** out of the distribution

Watch ...

Which of the five choices do we choose

Randomly pick one based on its distribution

To do this we need to **normalize** **<start>** out of the distribution

Watch ...

	i	to	chinese	food	lunch
<start>	0.015	0.01	0.005	0.003	0.001

$$0.015 + 0.01 + 0.005 + 0.003 + 0.001 = 0.034$$

<u>0.015</u>	<u>0.01</u>	<u>0.005</u>	<u>0.003</u>	<u>0.001</u>	▶
0.034	0.034	0.034	0.034	0.034	
0.44	0.29	0.15	0.09	0.03	

Now

We are going to use this sum to **normalize** out **<start>**

This is a divide sign

Which of the five choices do we choose

Randomly pick one based on its distribution

	i	to	chinese	food	lunch
<start>	0.015	0.01	0.005	0.003	0.001

$$0.015 + 0.01 + 0.005 + 0.003 + 0.001 = 0.034$$

To do this we need to **normalize** <start> out of the distribution

<u>0.015</u>	<u>0.01</u>	<u>0.005</u>	<u>0.003</u>	<u>0.001</u>
0.034	0.034	0.034	0.034	0.034
0.44	0.29	0.15	0.09	0.03

Now
our probabilities
sum to 1

Watch ...

Which of the **five** choices do we choose

	i	to	chinese	food	lunch
<start>	0.015	0.01	0.005	0.003	0.001

$$0.015 + 0.01 + 0.005 + 0.003 + 0.001 = 0.034$$

<u>0.015</u>	<u>0.01</u>	<u>0.005</u>	<u>0.003</u>	<u>0.001</u>
0.034	0.034	0.034	0.034	0.034
0.44	0.29	0.15	0.09	0.03

These are the probabilities you would pull the words from the bucket



Given a bucket with only those words that come after <start> we have probability 0.03 of picking *lunch*.

Which of the five choices do we choose

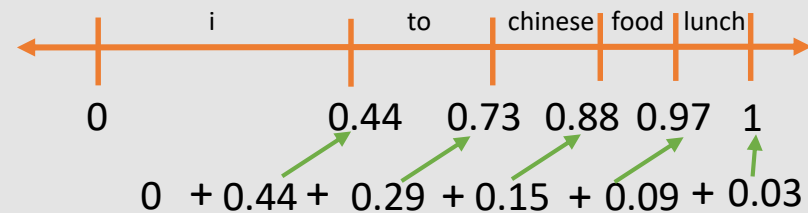
	i	to	chinese	food	lunch
<start>	0.015	0.01	0.005	0.003	0.001

$$0.015 + 0.01 + 0.005 + 0.003 + 0.001 = 0.034$$

<u>0.015</u>	<u>0.01</u>	<u>0.005</u>	<u>0.003</u>	<u>0.001</u>
0.034	0.034	0.034	0.034	0.034
0.44	0.29	0.15	0.09	0.03

	i	to	chinese	food	lunch
<start>	0.44	0.29	0.15	0.09	0.03

We can plot these probabilities
a line from 0 to 1



Which of the five choices do we choose

To pick a word that follows <start>:

Pick a **random number**
between 0 and 1,
and then see where it falls on the
distribution.

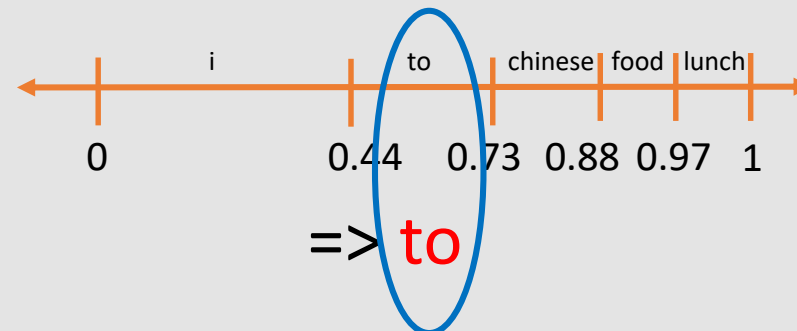
So say our random number
generator returned the value **0.65**,
what is our next word?

	i	to	chinese	food	lunch
<start>	0.015	0.01	0.005	0.003	0.001

$$0.015 + 0.01 + 0.005 + 0.003 + 0.001 = 0.034$$

<u>0.015</u>	<u>0.01</u>	<u>0.005</u>	<u>0.003</u>	<u>0.001</u>
0.034	0.034	0.034	0.034	0.034
0.44	0.29	0.15	0.09	0.03

	i	to	chinese	food	lunch
<start>	0.44	0.29	0.15	0.09	0.03



So then we start the process **again** with 'to'

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	0.002	0.33	0	0.0036	0	0	0	0.00079	0	0
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011	0	0
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087	0	0
eat	0	0	0.0027	0	0.021	0.0027	0.056	0	0	0.011
chinese	0.0063	0	0	0	0	0.52	0.0063	0	0	0.008
food	0.014	0	0.014	0	0.00092	0.0037	0	0	0	0.004
lunch	0.0059	0	0	0	0	0.0029	0	0	0	0.003
spend	0.0036	0	0.0036	0	0	0	0	0	1	0.006
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0
<end>	0	0	0	0	0.001	0.007	0.002	0.011	0	0



How to Generate Sentences?

The Shannon Visualization Method

• Four Steps

- 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**

```

<s> I
    I want
      want to
        to eat
          eat Chinese
            Chinese food
              food </s>

I want to eat Chinese food
  
```




Approximating Shakespeare

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'd in;
 Will you not tell me who I am?
 It cannot be but so.
 Indeed the short and the long. Marry, 'tis a noble Lepidus.



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)
- Quadrigrams worse: What's coming out looks like Shakespeare because it *is* Shakespeare

Zeros & Sparsity



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



Zeros

- Training set:
 - ... denied the allegations
 - ... denied the reports
 - ... denied the claims
 - ... denied the request
- Test set
 - ... denied the offer
 - ... denied the loan

$$P(\text{"offer"} \mid \text{denied the}) = 0$$

Sparsity

As **N** increases, the **accuracy** of our model increases

But

As **N** increases, the **sparsity** of our model increases

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	0.002	0.33	0	0.0036	0	0	0	0.00079	0	0
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011	0	0
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087	0	0
eat	0	0	0.0027	0	0.021	0.0027	0.056	0	0	0.011
chinese	0.0063	0	0	0	0	0.52	0.0063	0	0	0.008
food	0.014	0	0.014	0	0.00092	0.0037	0	0	0	0.004
lunch	0.0059	0	0	0	0	0.0029	0	0	0	0.003
spend	0.0036	0	0.0036	0	0	0	0	0	1	0.006
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0
<end>	0	0	0	0	0.001	0.007	0.002	0.011	0	0

LOOK AT ALL THE **ZEROS**

Does this mean that $P(\text{want} | \text{english}) = 0$?

With the model, yes but **in real life?**

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	0.002	0.33	0	0.0036	0	0	0	0.00079	0	0
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0064	0.0011	0	0
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087	0	0
eat	0	0	0.0027	0	0.021	0.0027	0.056	0	0	0.011
chinese	0.0063	0	0	0	0	0.52	0.0063	0	0	0.008
food	0.014	0	0.014	0	0.00092	0.0037	0	0	0	0.004
lunch	0.0059	0	0	0	0	0.0029	0	0	0	0.003
spend	0.0036	0	0.0036	0	0	0	0	0	1	0.006
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0
<end>	0	0	0	0	0.001	0.007	0.002	0.011	0	0

$$\begin{aligned}
 P(I \text{ want to eat English Food}) &= \\
 P(i | <start>) &* \\
 P(want | i) &* \\
 P(to | want) &* \\
 P(eat | to) &* \\
 P(english | eat) &* \\
 P(food | english) &* \\
 P(<end> | food) &= ?
 \end{aligned}$$

$$\begin{aligned}
 P(i | <start>) &= 0.015 \\
 P(want | i) &= 0.33 \\
 P(to | want) &= 0.66 \\
 P(eat | to) &= 0.28 \\
 P(english | eat) &= 0 \\
 P(food | english) &= 0 \\
 P(<end> | food) &= 0.004
 \end{aligned}$$

Uh oh!

$$P(I \text{ want to eat English Food}) = 0 ?$$

Sparsity & Recap MLE

- **Sparsity** is a major problem for **Maximum Likelihood Estimation (MLE)**
- This is **MLE** $\Rightarrow P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$
- MLE with example:

$$P(\text{the magical unicorn}) =$$

$$P(\text{the}) * \\ P(\text{magical} | \text{the}) * \\ P(\text{unicorn} | \text{magical})$$

These probabilities are referred to as **Relative Frequency**

Relative Frequency

- **Sparsity** is a major problem for **Maximum Likelihood Estimation (MLE)**
- This is **MLE** $\Rightarrow P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$
-

This is **Relative Frequency** $\Rightarrow P(w_k | w_{k-1})$

$$P(\text{unicorn} | \text{magical}) = \frac{\text{Frequency}(\text{magical unicorn})}{\text{Frequency}(\text{magical})}$$

Sparsity

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	0.002	0.33	0	0.0036	0	0	0	0.00079	0	0
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0064	0.0011	0	0
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087	0	0
eat	0	0	0.0027	0	0.021	0.0027	0.056	0	0	0.011
chinese	0.0063	0	0	0	0	0.52	0.0063	0	0	0.008
food	0.014	0	0.014	0	0.00092	0.0037	0	0	0	0.004
lunch	0.0059	0	0	0	0	0.0029	0	0	0	0.003
spend	0.0036	0	0.0036	0	0	0	0	0	1	0.006
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0
<end>	0	0	0	0	0.001	0.007	0.002	0.011	0	0

Because we don't see "<start> eat" in the text does this mean it doesn't occur ever?

Is $P(\text{<start> eat})$ really zero?

N-gram Smoothing

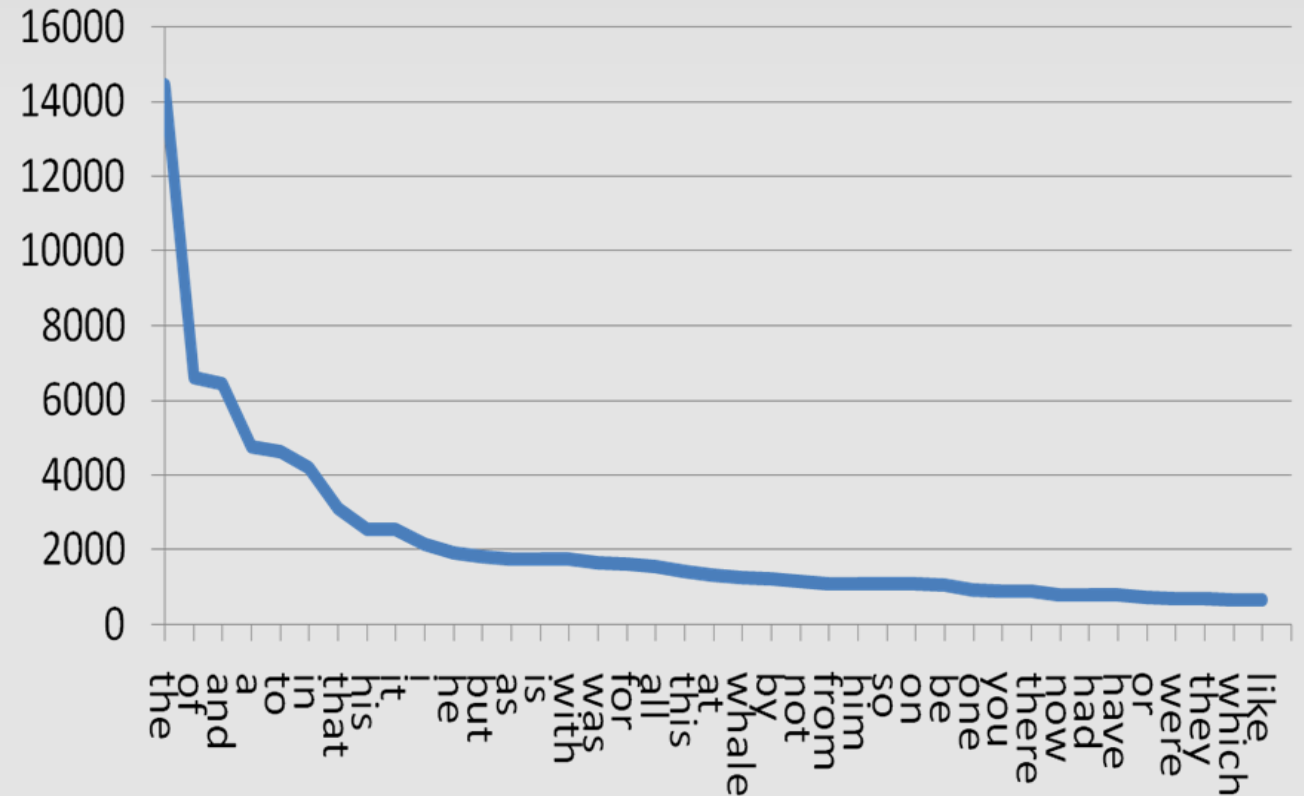
Smoothing

- Exploit the **Zipfian distribution** of words
- Two smoothing methods:
 - Laplace Smoothing
 - Good-Turning
- **The basic idea** is that we take a little from everything we see and give it to what we don't see
- Robin Hood: **stealing from the rich and giving to the poor**



Zipfian Distribution

- Words follow a Zipfian Distribution
 - Small number of words occur very frequently
 - A large number of words are only seen once
- Zipf's Law
 - A word's frequency is approximately inversely proportional to its rank in the word distribution list
- Great Video on Zipf's Law
 - <https://www.youtube.com/watch?v=fCn8zs912OE>



Laplace Smoothing (Add-1 smoothing)

- Simple metric : adds **one** to each count
- Pretend we saw each word one more time than we did

$$P(w_i) = \frac{\text{frequency}(w_i)}{N} \quad P_{\text{Laplace}}(w_i) = \frac{\text{frequency}(w_i) + 1}{N + \textcolor{blue}{V}}$$

N = the number of **tokens** in our corpus
V = the number of **types** in our corpus

Adding V because you've added one to each w seen in your corpus

Discounted Frequencies

$$P_{Laplace}(w_i) = \frac{\textit{frequency}(w_i) + \mathbf{1}}{N + \mathbf{V}}$$

$$P_{Laplace}(w_i) = \frac{\textit{frequency}^*(w_i)}{N}$$

$$\textit{frequency}^*(w_i) = (\textit{frequency}(w_i) + \mathbf{1}) \frac{N}{N + \mathbf{V}}$$

N = the number of **tokens** in our corpus

V = the number of **types** in our corpus

Discounted Frequencies and Probabilities

$$\text{frequency}^*(w_i) = (\text{frequency}(w_i) + 1) \frac{N}{N + V}$$

$$P_{\text{Laplace}}(w_i) = \frac{\text{frequency}^*(w_i)}{N}$$

$$1: P_{\text{Laplace}}(w_i) = \frac{(\text{frequency}(w_i) + 1) \frac{N}{N + V}}{N}$$

$$3: P_{\text{Laplace}}(w_i) = \frac{N(\text{frequency}(w_i) + 1)}{N + V} * \frac{1}{N}$$

$$2: P_{\text{Laplace}}(w_i) = \frac{\frac{N(\text{frequency}(w_i) + 1)}{N + V}}{N}$$

$$4: P_{\text{Laplace}}(w_i) = \frac{(\text{frequency}(w_i) + 1)}{N + V}$$

Laplace Smoothing on Conditional Probabilities

$$P(w_i) = \frac{\text{frequency}(w_i)}{N} \Rightarrow P_{\text{Laplace}}(w_i) = \frac{\text{frequency}(w_i) + \mathbf{1}}{N + \mathbf{V}}$$

$$P(w_1|w_2) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)} \Rightarrow P_{\text{Laplace}}(w_i|w_{i-1}) = ?$$

$$P_{\text{Laplace}}(w_n|w_{n-1}) = \frac{\text{frequency}(w_{n-1} w_n) + \mathbf{1}}{\text{frequency}(w_{n-1}) + \mathbf{V}}$$

\mathbf{V} = the number of **types** in our corpus

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	5	827	0	9	0	0	0	2	0	0
want	2	0	608	1	6	6	5	1	0	0
to	2	0	4	686	2	0	6	211	0	0
eat	0	0	2	0	16	2	42	0	0	34
chinese	1	0	0	0	0	82	1	0	0	23
food	15	0	15	0	1	4	0	0	0	12
lunch	2	0	0	0	0	1	0	0	0	9
spend	1	0	1	0	0	0	0	0	1	17
<start>	45	0	30	0	15	10	3	0	0	0
<end>	0	0	0	0	3	23	6	34	0	0

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

V = 1446

$$P_{Laplace}(want|i) = \frac{frequency(i\ want) + 1}{frequency(i) + V}$$

$$= \frac{827+1}{2533+1446} = 0.21$$

i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
2533	927	2417	746	158	1093	341	278	3000	3000

	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

	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

How smoothing affects probabilities?

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	0.002	0.33	0	0.0036	0	0	0	0.00079	0	0
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0064	0.0011	0	0
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087	0	0
eat	0	0	0.0027	0	0.021	0.0027	0.056	0	0	0.011
chinese	0.0063	0	0	0	0	0.52	0.0063	0	0	0.008
food	0.014	0	0.014	0	0.00092	0.0037	0	0	0	0.004
lunch	0.0059	0	0	0	0	0.0029	0	0	0	0.003
spend	0.0036	0	0.0036	0	0	0	0	0	1	0.006
<start>	0.015	0	0.01	0	0.005	0.003	0.001	0	0	0
<end>	0	0	0	0	0.001	0.007	0.002	0.011	0	0



Add-1 Estimation is a Blunt Instrument

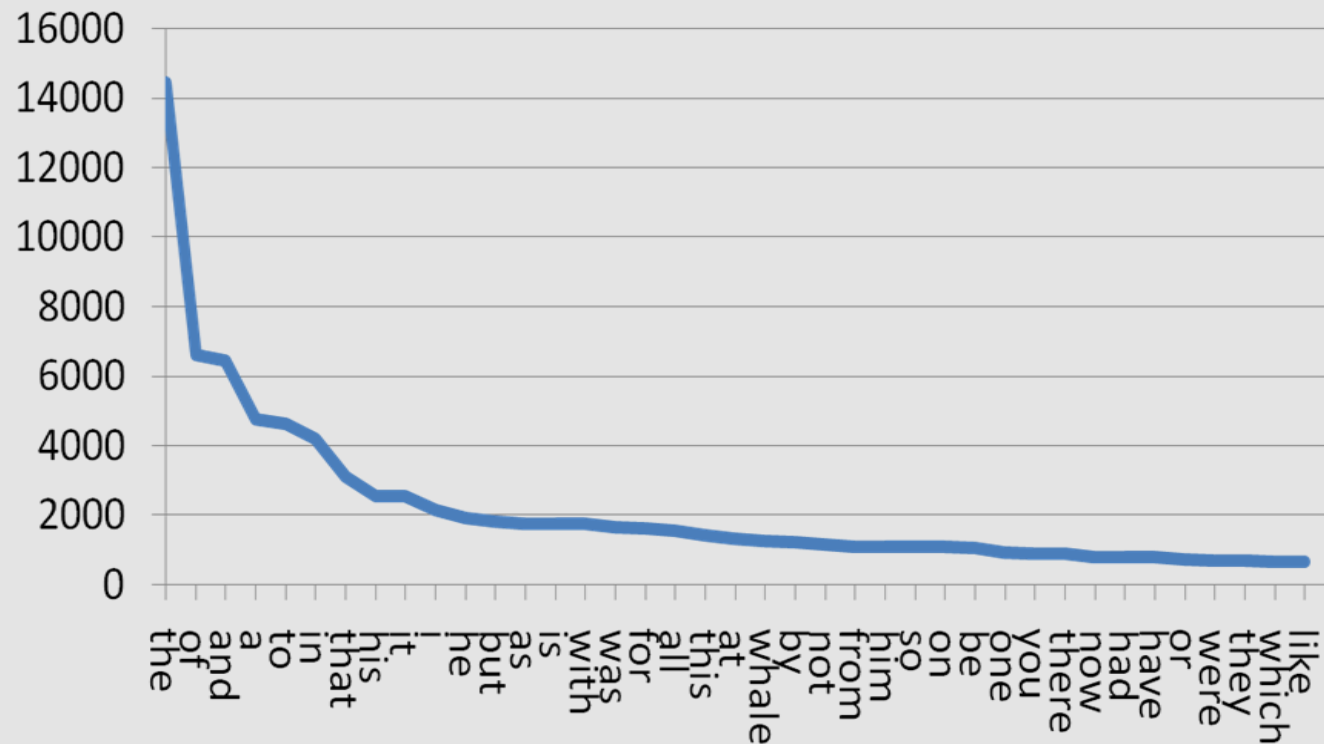
- So **add-1** isn't used for N-grams:
 - We'll see better methods
- But **add-1** is used to smooth other NLP models
 - For **text classification**
 - In domains where the **number of zeros isn't so huge.**

Good-Turing

- Add-1 smoothing (Laplace smoothing) is a bit **brute force**
- Few more **elegant** ways to smooth
 - **Good-Turning**
 - Witten-Bell
 - Kneser-Ney

Good-Turing

- Intuition
 - Use the count of things you have seen *once* to help estimate the count of things you've *never seen*



Good-Turing

Based on computing N_c which is the number of N-grams that occur c times

frequency of frequency


$N_0 = \# \text{ of bigrams with count } 0$

$N_1 = \# \text{ of bigrams with count } 1$


...

$N_c = \# \text{ of bigrams with count } c$

Adjust frequencies

$$\text{frequency}^*(w_i) = (\text{frequency}(w_i) + 1) \frac{N}{N + V}$$


Laplace Smoothing

$$\text{frequency}^*(w_i) = (\text{frequency}(w_i) + 1) \frac{N_{c+1}}{N_c}$$


Good Turing Smoothing

$$P_{\text{smoothing}}(w_n | w_{n-1}) = \frac{\text{frequency}^*(w_{n-1}w_n)}{\text{frequency}^*(w_{n-1})}$$

But what about **unseen** bigrams?

$$P_{gt}(\textit{unseen}) = \frac{N_1}{N_o}$$

N_1 = number of bigrams seen **1** time

N_o = **total** number of bigrams in the corpus

Example

Frequency	Frequency(Frequency)
0	2081496
1	5315
2	1419
3	642
4	381
5	311
6	196
Frequency	Frequency(Frequency)
...	...
2533	2
2534	2
....
M	1

Unigram

i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
2533	927	2417	746	158	1093	341	278	3000	3000

Bigram

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	5	827	0	9	0	0	0	2	0	0

i spend occurs **twice** in our corpus

$$\text{frequency}^*(w_i) = (\text{frequency}(w_i) + 1) \frac{N_{c+1}}{N_c}$$

$$\text{frequency}^*(i \text{ spend}) = (\text{frequency}(i \text{ spend}) + 1) \frac{N_3}{N_2}$$

$$\text{frequency}^*(i \text{ spend}) = (2 + 1) \frac{642}{1419} = 1.36$$

$$P_{gt}(\text{spend} | i) = \frac{\text{frequency}^*(i \text{ spend})}{\text{frequency}^*(i)} = \frac{1.36}{2534} = 0.00054$$

How do we know this?

$$\text{frequency}^*(i) = (\text{frequency}(i) + 1) \frac{N_{2534}}{N_{2533}} = (2533 + 1) \frac{2}{2} = 2534$$

Frequency	Frequency(Frequency)
0	2081496
1	5315
2	1419
3	642
4	381
5	311
6	196
Frequency	Frequency(Frequency)
...	...
2533	2
2534	2
....
M	1

	i	want	to	eat	chinese	food	lunch	spend	<start>	<end>
i	5	827	0	9	0	0	0	2	0	0

i spend occurs twice in our corpus

$$frequency^*(w_i) = (frequency(w_i) + 1) \frac{N_{c+1}}{N_c}$$

$$frequency^*(i\ spend) = (frequency(i\ spend) + 1) \frac{N_3}{N_2}$$

$$frequency^*(i\ spend) = (2 + 1) \frac{642}{1419} = 1.36$$

$$P_{gt}(spend | i) = \frac{frequency^*(i\ spend)}{frequency^*(i)} = \frac{1.36}{2534} = 0.00054$$

$$P^*(i\ to) = \frac{N_1}{N_0} = \frac{5315}{2081496} = 0.003$$



Unseen

What happens when $N_{c+1} = 0$

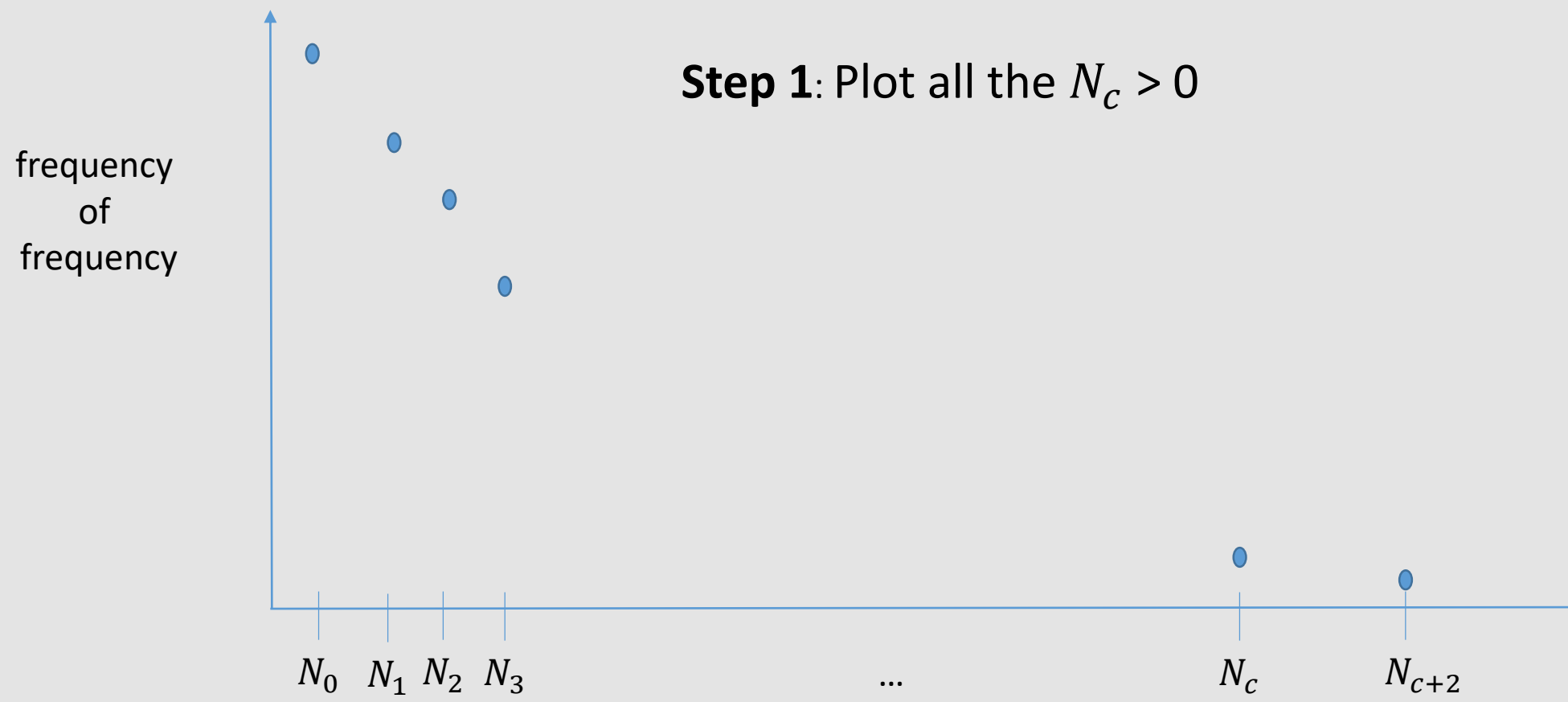
Frequency	Frequency(Frequency)
0	2081496
1	5315
2	1419
3	642
4	381
5	311
6	196
Frequency	Frequency(Frequency)
...	...
2533	2
2534	2
2535	0
....

$$frequency^*(w_i) = (frequency(w_i) + 1) \frac{N_{c+1}}{N_c}$$

Simplest thing is to perform **linear regressions** and replace the value of N_{c+1} with regression value whenever $N_{c+1} = 0$

Frequency of frequency
Is the number of n-grams
That occurred N_{c+1} times

Estimating when $N_{c+1} = 0$

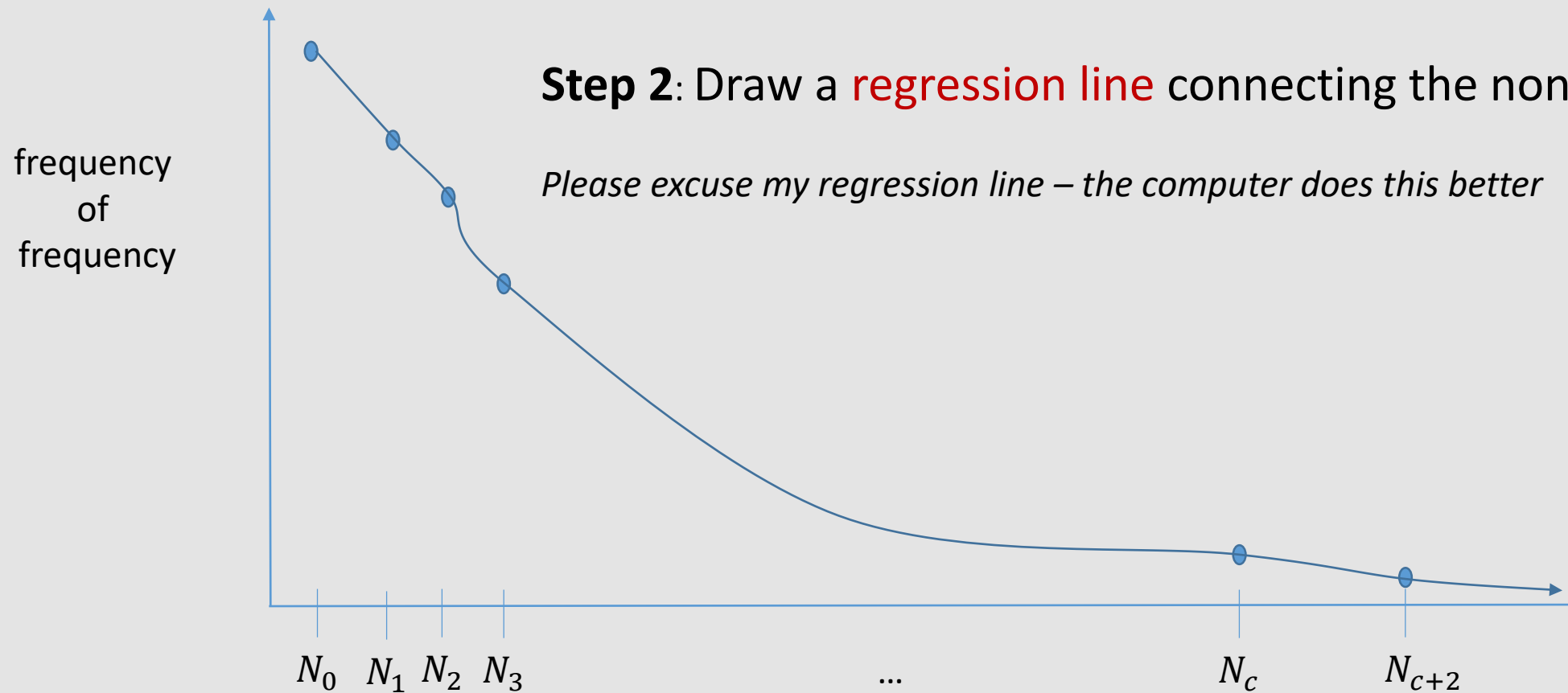


Frequency of frequency
Is the number of n-grams
That occurred N_{c+1} times

Estimating when $N_{c+1} = 0$

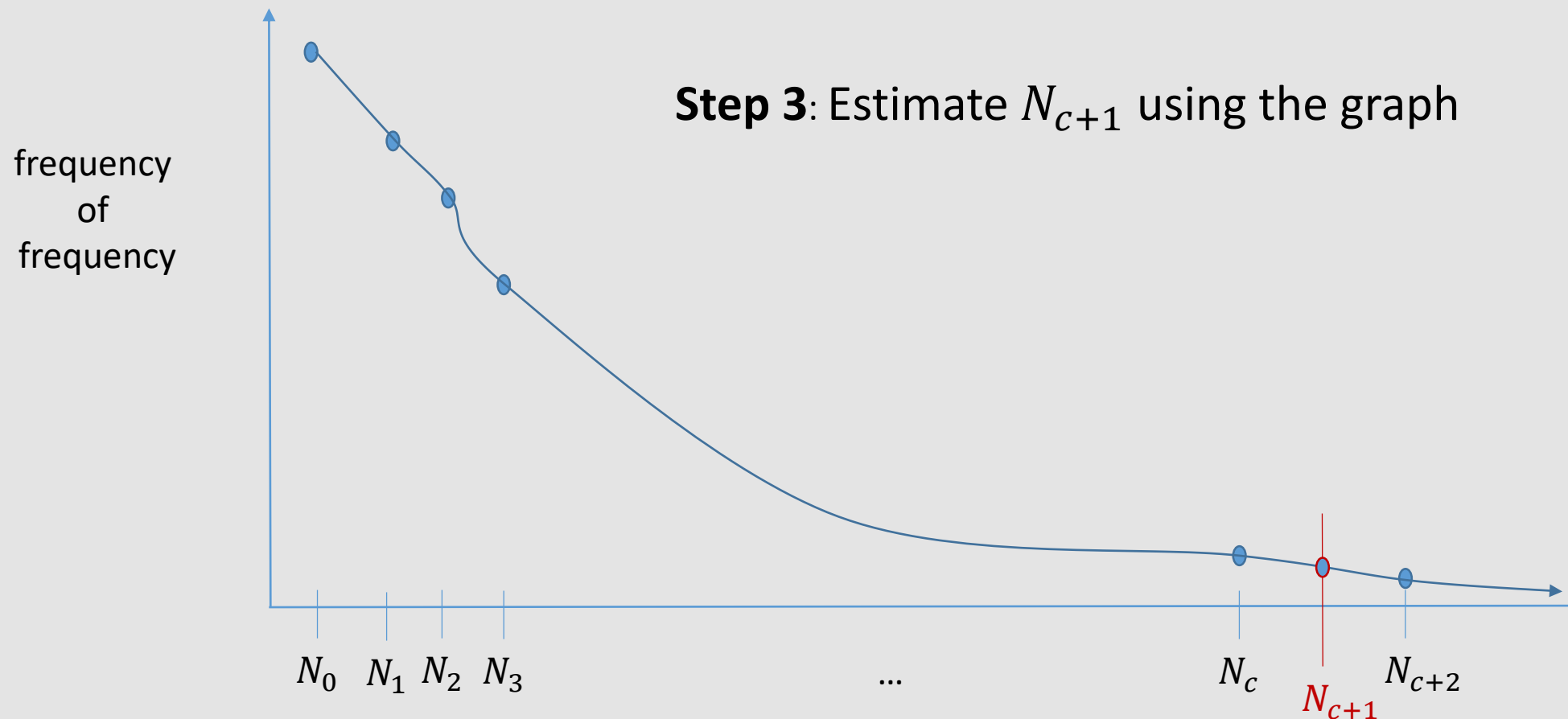
Step 2: Draw a **regression line** connecting the non-zero points

Please excuse my regression line – the computer does this better



Frequency of frequency
Is the number of n-grams
That occurred N_{c+1} times

Estimating when $N_{c+1} = 0$

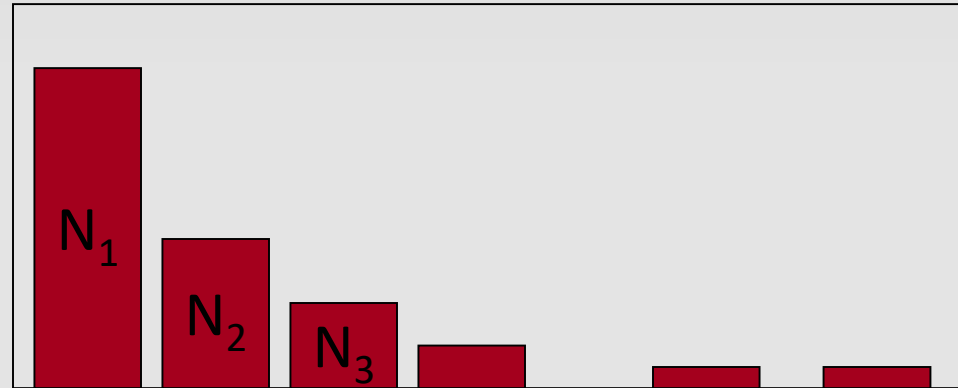


Good-Turing Complications

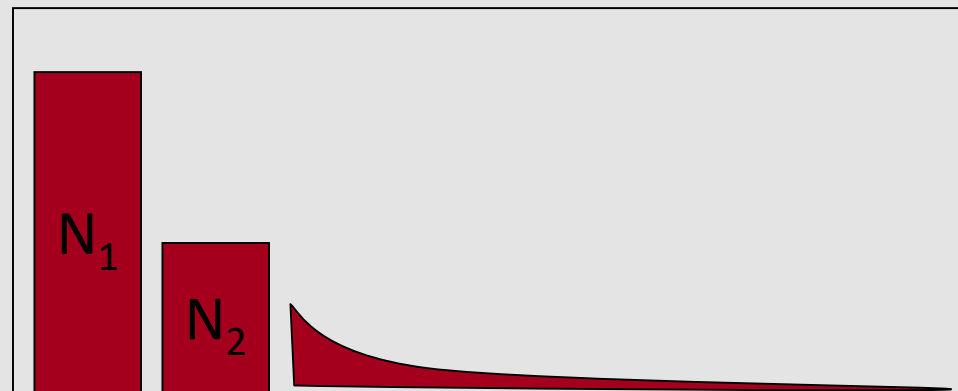
(slide from Dan Klein)

- Problem: what about “the”?
(say $c=4417$)

- For small k , $N_k > N_{k+1}$
- For large k , too jumpy, zeros wreck estimates



- Simple Good-Turing [Gale and Sampson]: replace empirical N_k with a best-fit power law once counts get unreliable





Huge Web-Scale N-grams

- How to deal with, e.g., Google N-gram corpus
- **Pruning**
 - **Only** store N-grams with **count > threshold**.
 - Remove singletons of higher-order n-grams
 - **Entropy-based** pruning
- **Efficiency**
 - Efficient **data structures**
 - Bloom filters: approximate language models
 - Store words as **indexes**, not strings
 - Use **Huffman coding** to fit large numbers of words into two bytes
 - **Quantize** probabilities (4-8 bits instead of 8-byte float)



Smoothing for Web-scale N-grams

- “Stupid backoff” (Brants *et al.* 2007)
- No discounting, just use **relative frequencies**

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

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



N-gram Smoothing Summary

- Add-1 smoothing:
 - OK for **text categorization/classification**, not for language modeling
- The most commonly used method:
 - Extended Interpolated Kneser-Ney
- For **very large N-grams** like the Web:
 - Stupid Backoff
 - Works **well** in practice

Next up

- Next time:
 - POS tagging (read Chapter 5)
 - Student presentations