

Ngrams ... + Smoothing

Lecture # 5

7 February 2018

Recap of Ngram Modeling Example

Creating our Relative Frequency Tables

CORPUS

Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl, but this weekend, the former NFL standout won't.

.....

**We are going to look at creating
a bigram model**

CORPUS

Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl, but this weekend, the former NFL standout won't.

.....

Unigram	Frequency
aaron	1
hernandez	1
helped	1
lead	1
the	3
new	1
england	1
patriots	1
into	1
2011	1
super	1
bowl	1
but	1
this	1
weekend	1
Former	1
NFL	1
standout	1
wo	1
n't	1
.	1

Unigram RAW FREQUENCIES

	Aaron	hernandez	helped	lead	the	new	england	patriots	into	2011	super	bowl	but	this	weekend	former	nfl	standout	wo	n't	.
aaron	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
hernandez	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
helped	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
lead	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
the	0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0
new	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
england	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
patriots	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
into	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2011	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
super	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
bowl	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
but	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
this	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
weekend	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Former	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
NFL	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
standout	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
wo	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
n't	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
.	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Bigram RAW FREQUENCIES

$$P(w_2 \mid w_1) = \text{Freq}(w_1 w_2) / \text{Freq}(w_1)$$

	aaron	hernandez	helped	lead	the	new	england	patriots	into	2011	super	bowl	but	this	weekend	former	nfl	standout	wo	n't	.
aaron	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
hernandez	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
helped	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
lead	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
the	0	0	0	0	0	1/3	0	0	0	1/3	0	0	0	0	0	1	0	0	0	0	0
new	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
england	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
patriots	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
into	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2011	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
super	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
bowl	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
but	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
this	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
weekend	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Former	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
NFL	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
standout	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
wo	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
n't	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
.	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

RELATIVE FREQUENCIES

P(Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl but this weekend the former NFL standout won't) =

$P(\text{aaron}) *$

$P(\text{hernandez} | \text{aaron}) *$

$P(\text{helped} | \text{hernandez}) *$

$P(\text{lead} | \text{helped}) *$

$P(\text{the} | \text{lead}) *$

$P(\text{new} | \text{the}) *$

$P(\text{england} | \text{new}) *$

....

$P(\text{n't} | \text{wo})$

P(Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl but this weekend the former NFL standout won't) =

$P(\text{aaron}) *$
 $P(\text{hernandez} | \text{aaron}) *$
 $P(\text{helped} | \text{hernandez}) *$
 $P(\text{lead} | \text{helped}) *$
 $P(\text{the} | \text{lead}) *$
 $P(\text{new} | \text{the}) *$
 $P(\text{england} | \text{new}) *$
....
 $P(\text{n't} | \text{wo})$

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

Chain rule of probability

$$P(X_1 \dots X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1^2) \dots P(X_n|X_1^{n-1})$$

P(Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl
but this weekend the former NFL standout won't) =

P(aaron) *

P(hernandez | aaron) *

P(helped | aaron hernandez) *

P(lead | aaron hernandez helped) *

P(the | aaron hernandez helped lead) *

....

P(n't | aaron hernandez ... standout wo)

Chain rule of probability

$$P(X_1 \dots X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1^2) \dots P(X_n|X_1^{n-1})$$

P(Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl
but this weekend the former NFL standout won't) =

P(aaron) *

P(hernandez | aaron) *

P(helped | aaron hernandez) *

P(lead | aaron hernandez helped) *

P(the | aaron hernandez helped lead) *

....

P(n't | aaron hernandez ... wo)

calculating this is
tough due to sparseness:

The chances of seeing "Aaron Hernandez helped lead the New England Patriots
into the 2011 Super Bowl but this weekend the former NFL standout won't" is slim

Markov assumption

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

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

Markov assumption

Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl but this weekend the former NFL standout won't

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

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

$P(n't | \text{aaron hernandez ...standout wo}) = P(n't | \text{wo})$ **Using a bigram model**

$P(n't | \text{aaron hernandez ...standout wo}) = P(n't | \text{standout wo})$ **Using a trigram model**

$P(n't | \text{aaron hernandez ...standout wo}) = P(n't | \text{nfl standout wo})$ **Using a 4-gram model**

etc ...

P(Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl, but this weekend, the former NFL standout) =

$P(\text{aaron})^*$

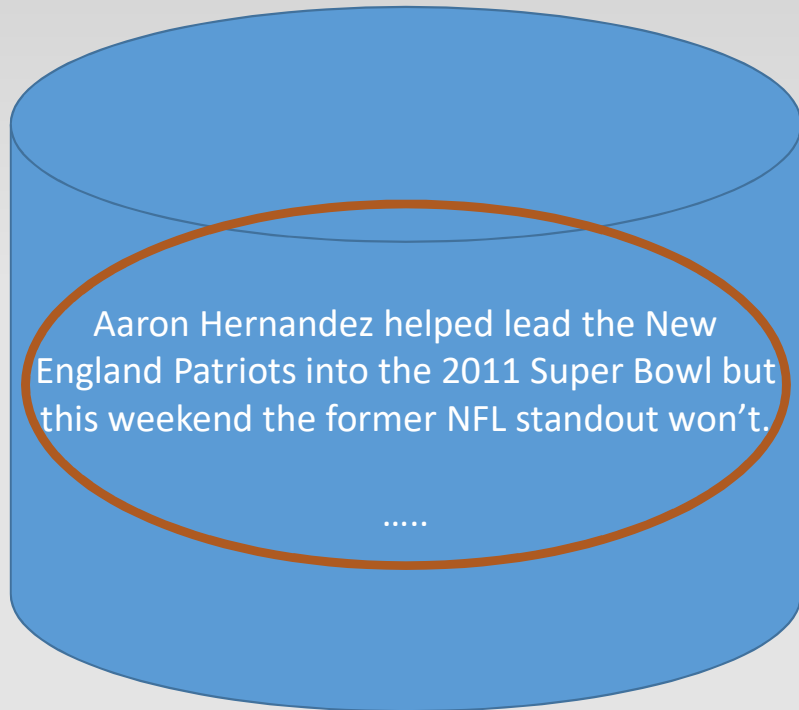
~~P(hernandez | aaron) *~~
$$P(\text{helped} \mid \text{Hernandez}) *$$
$$P(\text{lead} \mid \text{helped}) *$$
$$P(\text{the} \mid \text{lead}) *$$
$$P(\text{new} \mid \text{the}) *$$
$$P(\text{england} \mid \text{new}) *$$

• • • •

$$P(n't \mid wo)$$

Bigram Relative Frequency Table

[illegible]



My
Tokenization

Unigram	Frequency
aaron	1
hernandez	1
helped	1
lead	1
the	3
new	1
england	1
patriots	1
into	1
2011	1
super	1
bowl	1
but	1
this	1
weekend	1
Former	1
nfl	1
standou	1
Wo	1
n't	1
.	1

N = 23

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

N = 23

Unigram	Frequency
aaron	1
hernandez	1
helped	1
lead	1
the	3
new	1
england	1
patriots	1
into	1
2011	1
super	1
bowl	1
but	1
this	1
weekend	1
Former	1
nfl	1
standout	1
won	1
't	1
.	1



Unigram	P(unigram)
aaron	1/23 = 0.04
hernandez	1/23 = 0.04
helped	1/23 = 0.04
lead	1/23 = 0.04
the	3/23 = 0.13
new	1/23 = 0.04
england	1/23 = 0.04
patriots	1/23 = 0.04
into	1/23 = 0.04
2011	1/23 = 0.04
super	1/23 = 0.04
bowl	1/23 = 0.04
but	1/23 = 0.04
this	1/23 = 0.04
weekend	1/23 = 0.04
Former	1/23 = 0.04
nfl	1/23 = 0.04
standout	1/23 = 0.04
won	1/23 = 0.04
't	1/23 = 0.04
.	1/23 = 0.04

P(Aaron Hernandez helped lead the New England Patriots into the 2011 Super Bowl, but this weekend, the former NFL standout) =

$P(\text{aaron}) *$
 $P(\text{Hernandez} | \text{aaron}) *$
 $P(\text{helped} | \text{Hernandez}) *$
 $P(\text{lead} | \text{helped}) *$
 $P(\text{the} | \text{lead}) *$
 $P(\text{new} | \text{the}) *$
 $P(\text{england} | \text{new}) *$
....
 $P(\text{standout} | \text{NFL})$

Get this from our relative frequency tables

Does that make sense?

- Now let's look at a different example and introduce
 - <start>
 - <end>

	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	1	0	0	0	12
lunch	2	0	0	0	0	0	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 frequency's

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

Unigram table of raw frequency's

	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	1	0	0	0	12
lunch	2	0	0	0	0	0	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

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

Relative Frequency

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$




$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

	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	1	0	0	0	12
lunch	2	0	0	0	0	0	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

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

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$



$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

$$P(\text{want} | i) = \frac{827}{2533} = 0.33$$

	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

Relative Frequency Table

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$



$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

For bigrams

How do we generalize this for all n – grams

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$



$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

For bigrams

Any n-gram

$$P(w_n|w_{n-N+1}^{n-1}) = \frac{\text{frequency}(w_{n-N+1}^{n-1}, w_n)}{\text{frequency}(w_{n-N+1}^{n-1})}$$

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$



$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

For bigrams

Any n-gram

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{\text{frequency}(w_{n-N+1}^{n-1}, w_n)}{\text{frequency}(w_{n-N+1}^{n-1})}$$

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$



$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

For bigrams

Any n-gram

$$P(w_n|w_{n-N+1}^{n-1}) = \frac{\text{frequency}(w_{n-N+1}^{n-1}, w_n)}{\text{frequency}(w_{n-N+1}^{n-1})}$$

Trigram Model: $P(\text{unicorns}|\text{the magical}) = \frac{\text{frequency}(\text{the magical unicorns})}{\text{frequency}(\text{the magical})}$

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$



$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

For bigrams

Any n-gram

$$P(w_n|w_{n-N+1}^{n-1}) = \frac{\text{frequency}(w_{n-N+1}^{n-1}, w_n)}{\text{frequency}(w_{n-N+1}^{n-1})}$$

Trigram Model: $P(\text{unicorns}|\text{the magical}) = \frac{\text{frequency}(\text{the magical unicorns})}{\text{frequency}(\text{the magical})}$

4-gram Model: $P(\text{unicorns}|\text{the heroic magical}) = \frac{\text{frequency}(\text{the heroic magical unicorns})}{\text{frequency}(\text{the heroic magical})}$

$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\sum_w \text{frequency}(w_1 w)}$$



$$P(w_2|w_1) = \frac{\text{frequency}(w_1 w_2)}{\text{frequency}(w_1)}$$

For bigrams

Any n-gram

$$P(w_n|w_{n-N+1}^{n-1}) = \frac{\text{frequency}(w_{n-N+1}^{n-1}, w_n)}{\text{frequency}(w_{n-N+1}^{n-1})}$$

Trigram Model: $P(\text{unicorns}|\text{the magical}) = \frac{\text{frequency}(\text{the magical unicorns})}{\text{frequency}(\text{the magical})}$

4-gram Model: $P(\text{unicorns}|\text{the heroic magical}) = \frac{\text{frequency}(\text{the heroic magical unicorns})}{\text{frequency}(\text{the heroic magical})}$

5-gram Model: $P(\text{unicorns}|\text{the heroic unseen magical}) = \frac{\text{frequency}(\text{the heroic unseen magical unicorns})}{\text{frequency}(\text{the heroic unseen magical})}$

To calculate these we use: two tables

N-gram table

and

(N-1)-gram tables

Any questions?

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	0.0064	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

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>
<start>	0.15	0	0.1	0	0.4	0.3	0.05	0	0	0

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>
<start>	0.15	0	0.1	0	0.4	0.3	0.05	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

Which of the five choices do we 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 its distribution

	i	to	chinese	food	lunch
<start>	0.15	0.1	0.4	0.3	0.05

Which of the five choices do we choose

Randomly pick one based on its distribution

	i	to	chinese	food	lunch
<start>	0.15	0.1	0.4	0.3	0.05

Which of the five choices do we choose

	i	to	chinese	food	lunch
<start>	0.15	0.1	0.4	0.3	0.05

$$0.15 + 0.1 + 0.4 + 0.3 + 0.05 = 1$$



We can plot these probabilities
a line from 0 to 1

Which of the five choices do we choose

	i	to	chinese	food	lunch
<start>	0.15	0.1	0.4	0.3	0.05

$$0.15 + 0.1 + 0.4 + 0.3 + 0.05 = 1$$



We can plot these probabilities
a line from 0 to 1

Now we can randomly pick what word
Follows <start> given the distribution
of them occurring within the text

Which of the five choices do we choose

Pick a random number between zero and one

```
my $r = rand();
```

And then see where it falls on the distribution:

```
if($r <= 0.44)      { $next_word = "i"; }  
elseif($r <= 0.73) { $next_word = "to"; }  
elseif($r <= 0.88){ $next_word = "chinese"; }  
elseif($r <= 0.97) { $next_word = "food"; }  
else   { $next_word = "lunch"; }
```

	i	to	chinese	food	lunch
<start>	0.15	0.1	0.4	0.3	0.05

$$0.15 + 0.1 + 0.4 + 0.3 + 0.05 = 1$$



We can plot these probabilities
a line from 0 to 1

Now we can randomly pick what word
Follows <start> given the distribution
of them occurring within the text

Which of the five choices do we choose

Pick a random number between zero and one

```
my $r = rand();
```

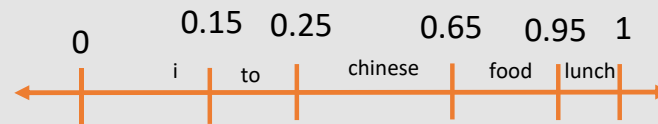
And then see where it falls on the distribution:

```
if($r <= 0.44)      { $next_word = "i"; }  
elseif($r <= 0.73) { $next_word = "to"; }  
elseif($r <= 0.88){ $next_word = "chinese"; }  
elseif($r <= 0.97) { $next_word = "food"; }  
else { $next_word = "lunch"; }
```

So say our rand returned the value 0.6 what is our next word?

	i	to	chinese	food	lunch
<start>	0.15	0.1	0.4	0.3	0.05

$$0.15 + 0.1 + 0.4 + 0.3 + 0.05 = 1$$



We can plot these probabilities
a line from 0 to 1

Now we can randomly pick what word
Follows <start> given the distribution
of them occurring within the text

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	0.0064	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

Do you see how this extends any n-gram?

	I am	to eat	chinese dish	...	lunch box
<start>	0.44	0.29	0.15	...	0.03

And before we move on

Perl's Hashes of Hashes

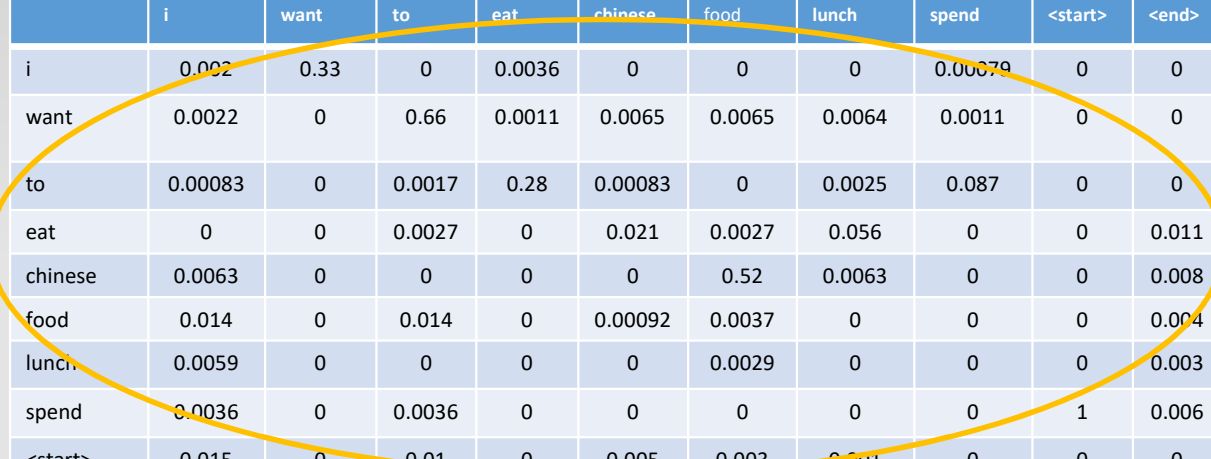
```
my $n = shift;
my %hash = ();
while(<>) {
    chomp;
    my @array = split/\s+/;
    for my $i(0..$#array) {
        my $j = $i + $n - 1;
        my $first = $array[$i];
        my $ngram = "";
        if($j > $#array) { next; }
        for my $k ($i..$j) {
            $ngram .= "$array[$k] ";
        }
        chomp $ngram;
        $hash{$first}{$ngram}++;
    }
}
```

Scarcity

As N increases the accuracy of our model increase

But

As N increases the sparsely 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	0.0064	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

LOOK AT ALL THE ZERO'S

	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

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

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

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

	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	0.0064	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

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

	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	0.0064	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

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

$$\begin{aligned}
 P(i | <start>) &= 0.015 \\
 P(want | i) &= 0.33 \\
 P(to | want) &= 0.66 \\
 P(eat | to) &= 0.28 \\
 P(english | eat) &= 0 \\
 P(food | english) &= 0 \\
 P(<end> | food) &= 0.007
 \end{aligned}$$

	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

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

$P(i | \text{<start>})$ *

$P(\text{want} | i)$ *

$P(\text{to} | \text{want})$ *

$P(\text{eat} | \text{to})$ *

$P(\text{english} | \text{eat})$ *

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

$P(\text{<end>} | \text{food}) = ?$

$P(i | \text{<start>}) = 0.015$

$P(\text{want} | i) = 0.33$

$P(\text{to} | \text{want}) = 0.66$

$P(\text{eat} | \text{to}) = 0.28$

$P(\text{english} | \text{eat}) = 0$

$P(\text{food} | \text{english}) = 0$

$P(\text{<end>} | \text{food}) = 0.007$

Uh oh!

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

.... ?

Sparicity

- Sparcity is a major problem for **Maximum Likelihood Estimation**

Recap MLE

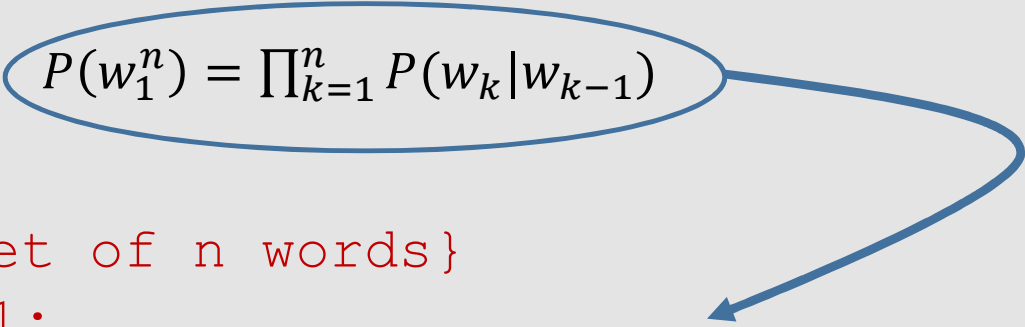
- Sparsity is a major problem for **Maximum Likelihood Estimation**

This is MLE =>
$$P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$$

MLE with code

- Sparsity is a major problem for **Maximum Likelihood Estimation**

This is MLE =>

$$P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$$


```
my @w = {set of n words}
my $P_w = 1;
for my $k (1...n) {
    $P_w *= $rel_freq_table{$w[$k]}{$w[$k-1]}
}
```

MLE with example

- Sparsity is a major problem for **Maximum Likelihood Estimation**

This is MLE => $P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$

$$\begin{aligned} P(\textit{the magical unicorn}) = \\ P(\textit{the}) * \\ P(\textit{magical} | \textit{the}) * \\ P(\textit{unicorn} | \textit{magical}) \end{aligned}$$

MLE with example

- Sparsity is a major problem for **Maximum Likelihood Estimation**

This is MLE => $P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$

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

$$P(\text{the}) *$$

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

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

These probabilities are referred to as Relative Frequency

Relative Frequency

- Sparsity is a major problem for **Maximum Likelihood Estimation**

This is MLE => $P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$

This is Relative Frequency => $P(w_k | w_{k-1})$

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

```
$rel_freq_table{$w[$k]}{$w[$k-1]}
```

Sparcity

	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

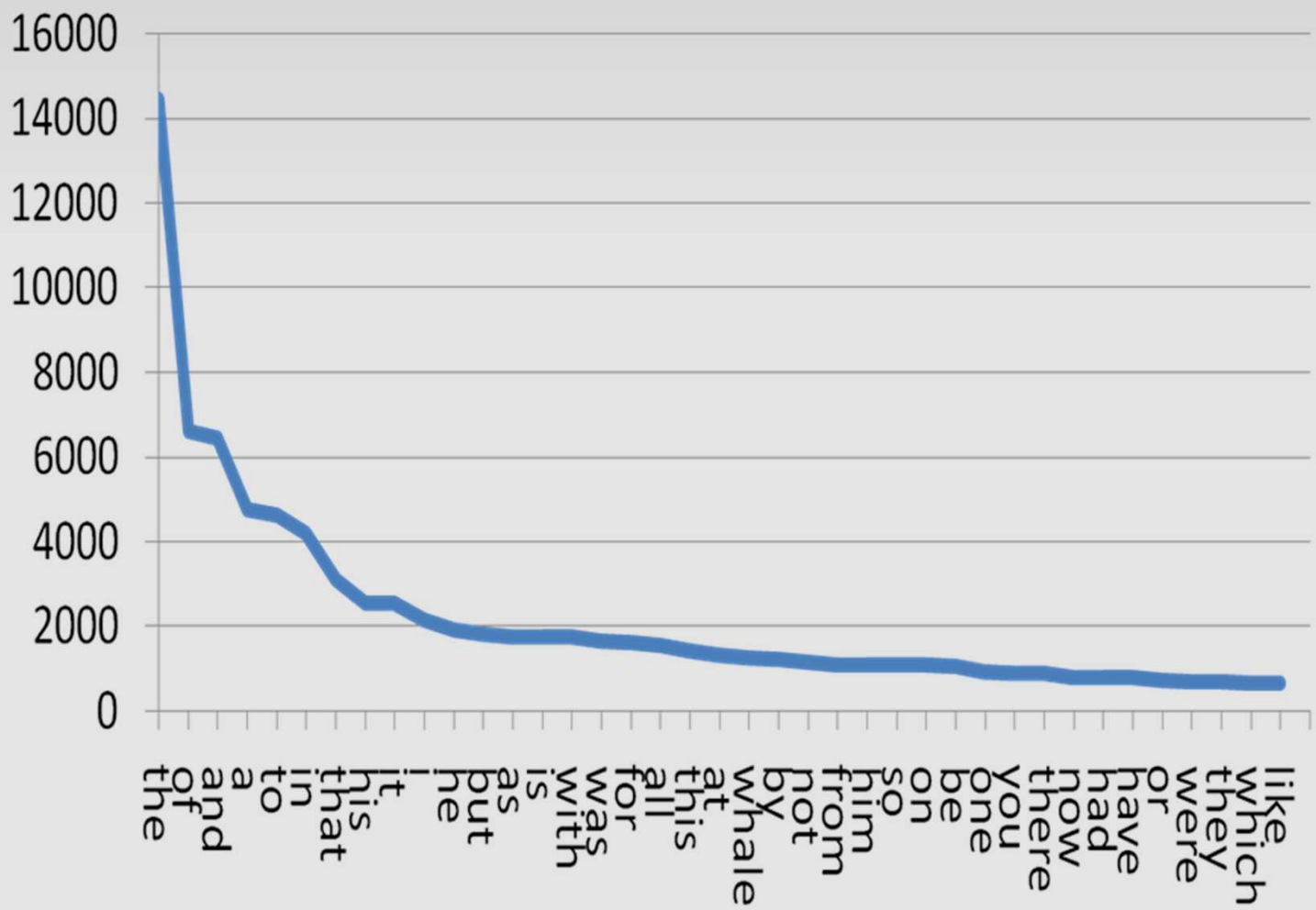
Smoothing

- 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

- <https://www.youtube.com/watch?v=fCn8zs912OE>

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



Laplace Smoothing

- Simple metric : adds one to each count

$$P(w_i) = \frac{\text{frequency}(w_i)}{N}$$

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

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

Laplace Smoothing

- Simple metric : adds one to each count

$$P(w_i) = \frac{\text{frequency}(w_i)}{N}$$

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

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

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

The book refers to “adjusted count”

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

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

Adjusted count

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

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

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

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

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

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

Adjusted count

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

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

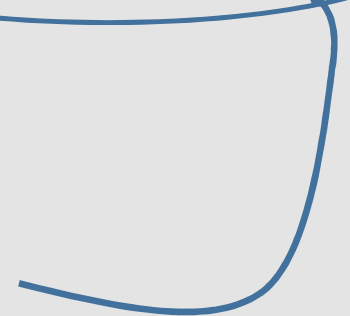
$$P_{\text{Laplace}}(w_i) = \frac{\text{frequency}(w_i) + 1}{N + V}$$

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

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

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

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



Laplace Smoothing on Conditional Probabilities

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

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

V = the number of types in our corpus

Laplace Smoothing on Conditional Probabilities

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

$$P(w_1|w_2) = \frac{\text{frequency}(w_1, w_2)}{\text{frequency}(w_2)} \Rightarrow P_{\text{Laplace}}(w_i) = ?$$

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

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	1	0	0	0	12
lunch	2	0	0	0	0	0	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

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

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

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

$$V = 1446$$

	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	1	0	0	0	12
lunch	2	0	0	0	0	0	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

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

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

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

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

$$V = 1446$$

	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	0.00056	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

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

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

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

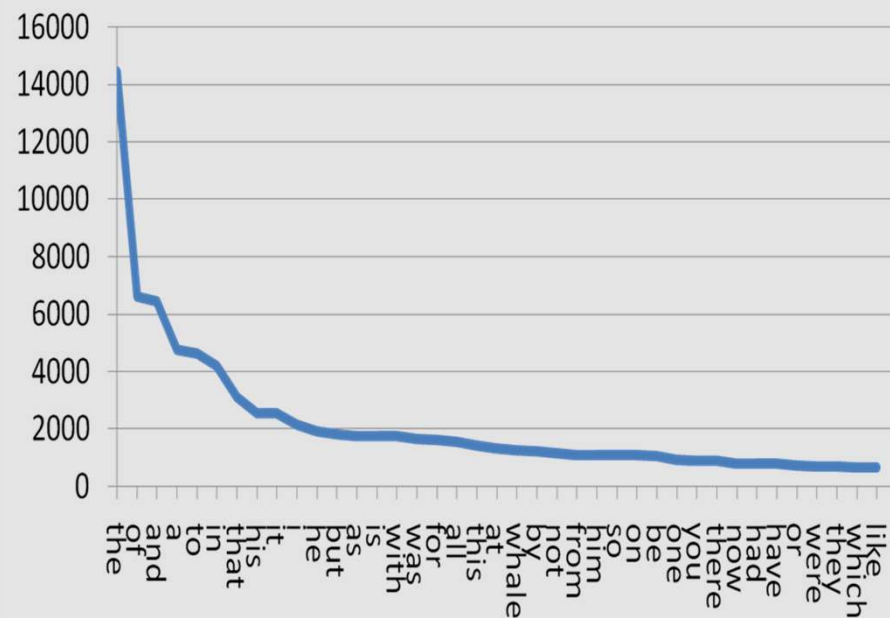
V = 1446

Good Turing

- Add-1 smoothing (Laplace smoothing) is a bit brute force
 - Few more elegant ways to smooth
 - **Good Turing**
 - 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

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

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

...

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

Redefine frequency

$$\text{frequency}^*(w_i) = (\text{frequency}(w_i) + 1) \frac{N}{N + V} \quad \left. \vphantom{\text{frequency}^*(w_i)} \right\} \text{Laplace Smoothing}$$

Redefine frequency

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

} Laplace Smoothing

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

} Good Turing Smoothing

Redefine frequency

$$\left. frequency^*(w_i) = (frequency(w_i) + 1) \frac{N}{N + V} \right\} \text{Laplace Smoothing}$$

$$\left. frequency^*(w_i) = (frequency(w_i) + 1) \frac{N_{c+1}}{N_c} \right\} \text{Good Turing Smoothing}$$

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

But what about unseen bigrams

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

How do we know what N_o is given
we don't know the number unseen events?

N_1 = number of bigrams seen 1 time
 N_o = total number of bigrams in the corpus

But what about unseen bigrams

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

How do we know what N is given
we don't know the number unseen events?

N_1 = number of bigrams seen 1 time
 N_o = total number of bigrams in the corpus

Guesstimate

We know V (the vocabulary size), therefore

The total number of bigrams = V^2

So

$$N_o = V^2 - \# \text{ seen bigrams}$$

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

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

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

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

$$P_{gt}(\text{spend} | i) = \frac{\text{frequency}^*(i \text{ spend})}{\text{frequency}^*(i)} = \frac{1.36}{2534} = 0.00054 \text{ (versus 0.00039)}$$

How do we know this?

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

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

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

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

$$P_{gt}(\text{spend} | i) = \frac{\text{frequency}^*(i \text{ spend})}{\text{frequency}^*(i)} = \frac{1.36}{2534} = 0.00054 \text{ (versus 0.00039)}$$

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

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

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

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

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

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

$$P_{gt}(\text{spend} | i) = \frac{\text{frequency}^*(i \text{ spend})}{\text{frequency}^*(i)} = \frac{1.36}{2534} = 0.00054 \text{ (versus } 0.00039)$$

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

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

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

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

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

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

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

$$P_{gt}(\text{spend} | i) = \frac{\text{frequency}^*(i \text{ spend})}{\text{frequency}^*(i)} = \frac{1.36}{2534} = 0.00054 \text{ (versus } 0.00039)$$

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

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

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

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

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

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

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

$$P_{gt}(\text{spend} | i) = \frac{\text{frequency}^*(i \text{ spend})}{\text{frequency}^*(i)} = \frac{1.36}{2534} = 0.00054 \text{ (versus } 0.00039)$$

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

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

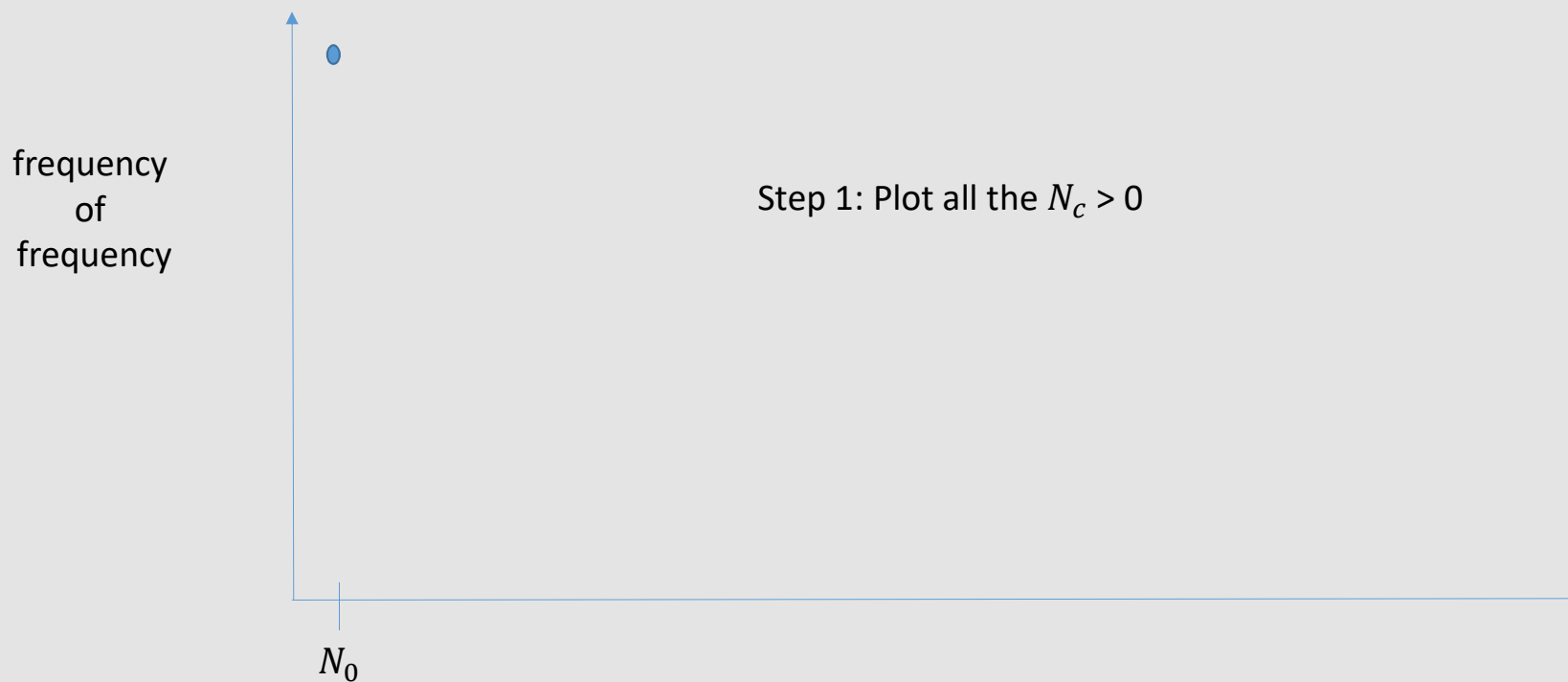
Frequency	Frequency(Frequency)
0	2081496
1	5315
2	1419
3	642
4	381
5	311
6	196
...
2533	2
2534	2
2535	0
....
M	1

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

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

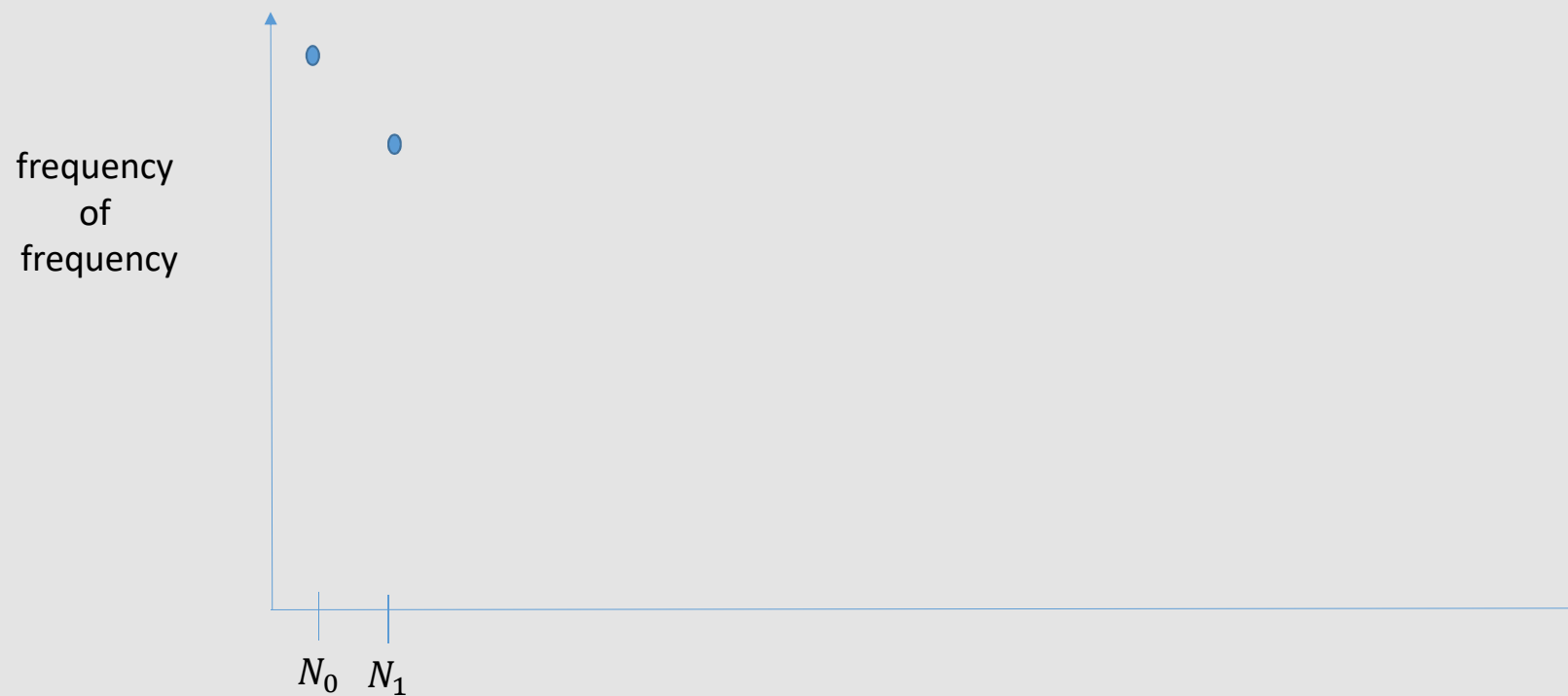
Frequency of frequency
Is the number of n-grams
That occurred N_{c+1} times

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



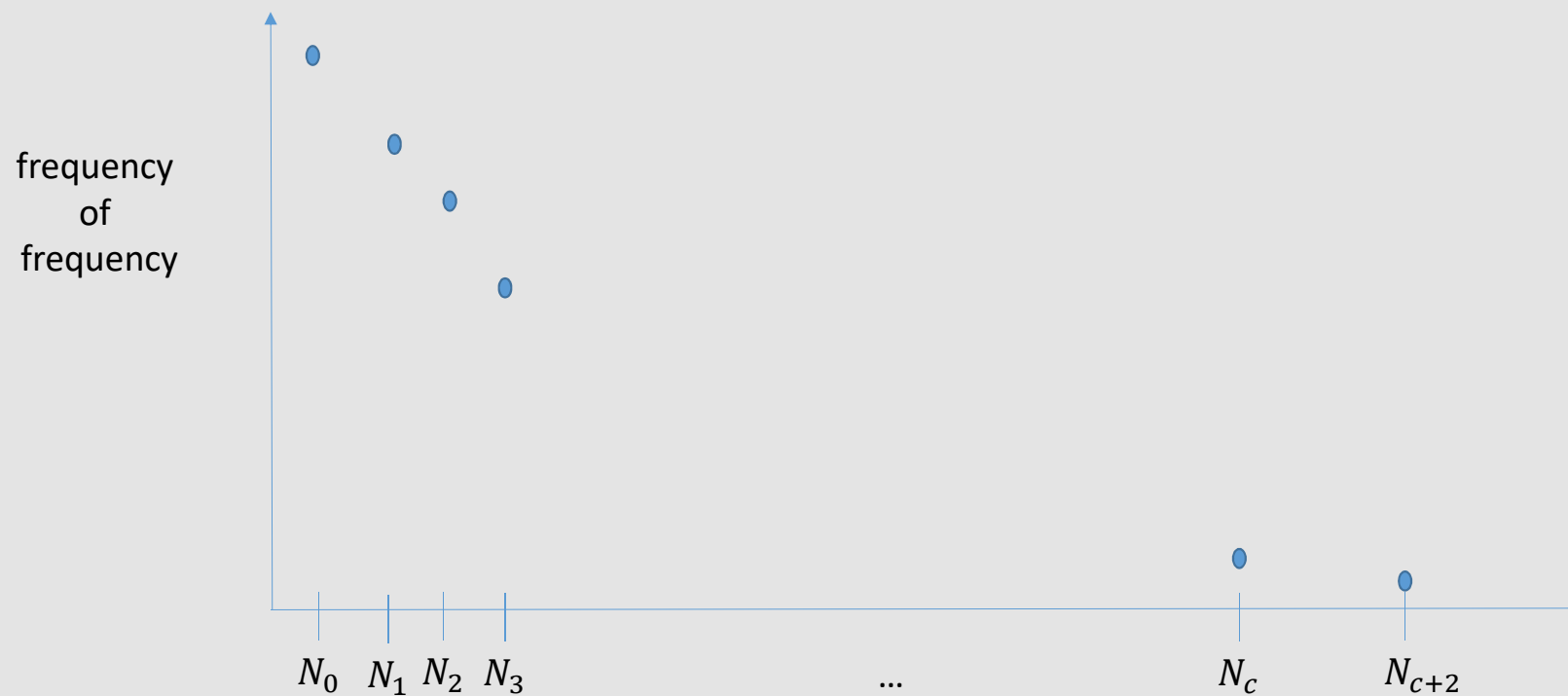
Frequency of frequency
Is the number of n-grams
That occurred N_{c+1} times

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



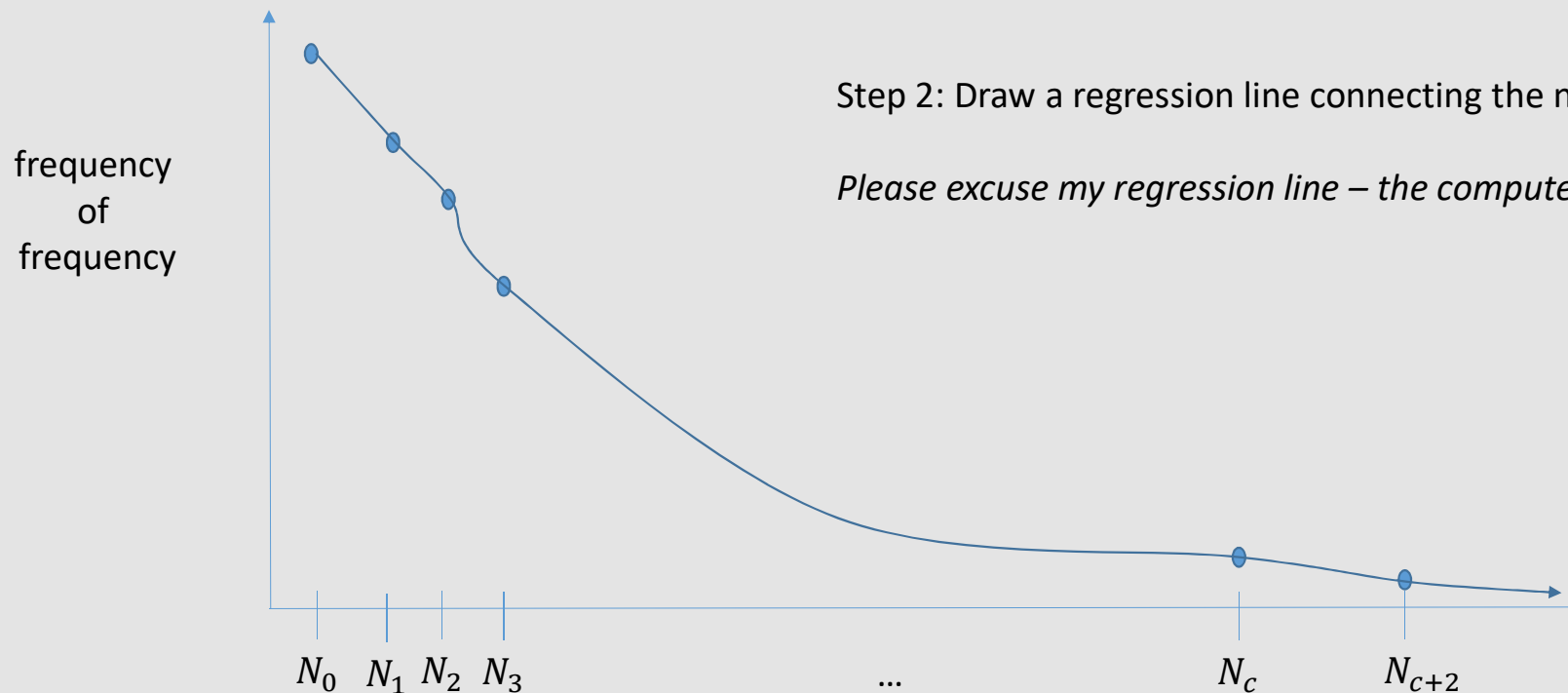
Frequency of frequency
Is the number of n-grams
That occurred N_{c+1} times

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



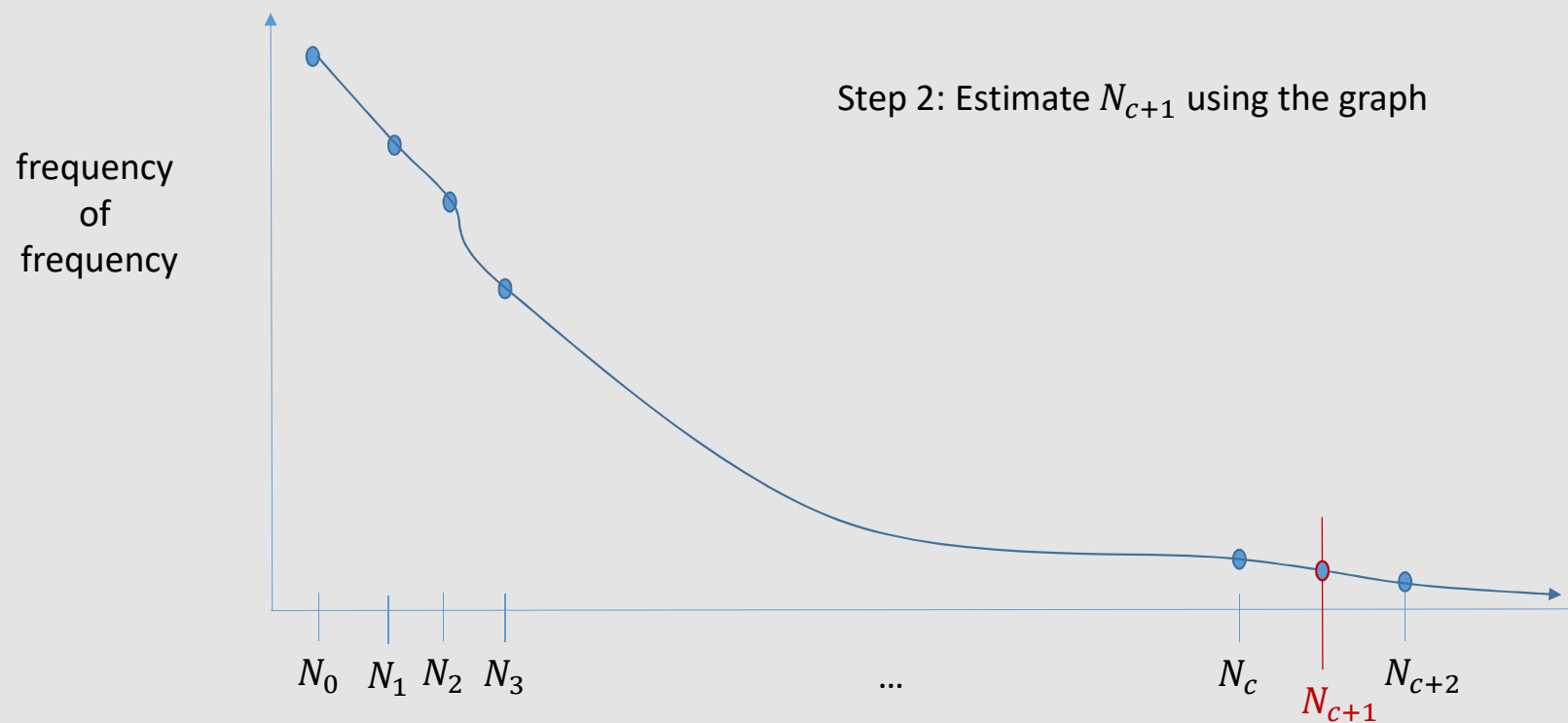
Frequency of frequency
Is the number of n-grams
That occurred N_{c+1} times

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



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

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



Questions?