

# N-gram Language Models

COMP90042

Natural Language Processing  
Lecture 3

Semester 1 2022 Week 2  
Jey Han Lau



THE UNIVERSITY OF  

---

MELBOURNE

# Language Models

- One NLP application is about *explaining language*
  - Why some sentences are more fluent than others
- E.g. in speech recognition:
  - *recognise speech > wreck a nice beach*
  - We measure ‘goodness’ using probabilities estimated by language models
  - Language model can also be used for generation
    - ↳ *N-gram*

**Generate Options**[Learn more in the docs.](#)Length to generate [?](#) 500 Start at beginning [?](#)[Advanced Settings »](#)

Donald Trump decides to become a vegan.

That's according to the diet of the President of the United States, who revealed his new eating habits Monday at a Mississippi campaign rally for Sen. Cindy Hyde-Smith, whom he endorsed over Democrat Mike Espy.

Trump, 71, spoke to the crowd at the Tupelo Regional Airport about his years of eating "perfectly healthy," including steaks and other meats.

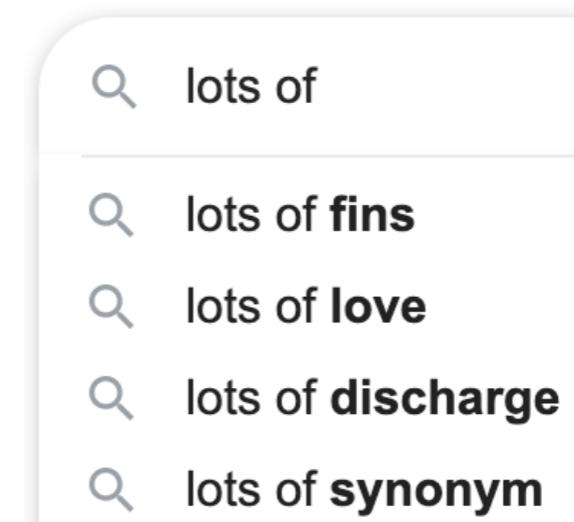
But now he's turned over a new leaf. "I'm not a fan of meat."

Trump told the crowd he now eats "mostly fish, like quite a bit of fish," and

[Generate Text](#)

# Language Models

- Useful for
  - Query completion
  - Optical character recog.
- Other generation tasks
  - Machine translation
  - Summarisation
  - Dialogue systems
- Nowadays pretrained language models are the backbone of modern NLP systems



# Outline

- Deriving  $n$ -gram language models
- Smoothing to deal with sparsity

# Probabilities: Joint to Conditional

Our goal is to get a probability for an arbitrary sequence of  $m$  words

$$P(w_1, w_2, \dots, w_m)$$

First step is to apply the chain rule to convert joint probabilities to conditional ones

$$P(w_1, w_2, \dots, w_m) = P(w_1)P(w_2 | w_1)P(w_3 | w_1, w_2) \dots$$
$$P(w_m | \underbrace{w_1, \dots, w_{m-1}}_{\text{red arrow}})$$

# The Markov Assumption

(stat est.)

Too hard to find the prob. of such 19 words in that order in corpus.  
i.e. 20 words in a sentence.  $P(w_0 | w_1 \dots w_{19})$

Still intractable, so make a simplifying assumption:

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | \underbrace{w_{i-n+1} \dots w_{i-1}}_{\text{just look the fixed prior context (i.e. 3)}}$$

For some small  $n$

When  $n = 1$ , a unigram model (bag-of-words model)

↓

compute each word independently

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i)$$

$w_i$

the dog barks

+ throw the sentence meaning away / just care some specific word  
When  $n = 2$ , a bigram model

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-1})$$

$w_i$

the dog barks

When  $n = 3$ , a trigram model

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-2} w_{i-1})$$

$w_i$

the dog barks

# Maximum Likelihood Estimation

How do we calculate the probabilities? Estimate based on counts in our corpus:

For unigram models,

How many times you see "barks" in the corpus  
 ↗

$$P(w_i) = \frac{C(w_i)}{M} \rightarrow \frac{C(\text{barks})}{M}$$

All word tokens in corpus

For bigram models,

How many times "dog barks" in yr corpus  
 ↗

$$P(w_i | w_{i-1}) = \frac{C(w_{i-1}w_i)}{C(w_{i-1})} \rightarrow \frac{C(\text{dog barks})}{C(\text{dog})}$$

↳ how many times "dog" showed up  
 (pre-condition)

~~★~~  $P(w_i | w_{i-n+1} \dots w_{i-1}) = \frac{C(w_{i-n+1} \dots w_i)}{C(w_{i-n+1} \dots w_{i-1})}$

<sup>20</sup> <sub>18, 19, 20</sub>

if n=3 . total word in sentence is 20

# Book-ending Sequences

- Special tags used to denote start and end of sequence
  - ▶  $<S>$  = sentence start
  - ▶  $</S>$  = sentence end

# Trigram example

Corpus:

*yes no no no yes  
no no no yes yes yes no*

What is the probability of

*< s > < s > yes no no yes < /s >*  
<sup>①</sup> <sub>④</sub>

under a trigram language model?

$P(\text{yes no no yes}) =$

$$P(\text{yes} | \langle s \rangle \langle s \rangle) \times$$

$$P(\text{no} | \langle s \rangle \text{ yes}) \times$$

$$P(\text{no} | \text{yes no}) \times$$

$$P(\text{yes} | \text{no no}) \times$$

$$P(\langle /s \rangle | \text{no yes})$$

Need to predict  $\langle /s \rangle$  because it's the end of sentence!

Corpus:

*<S> <S> yes no no no no yes </S>*  
*<S> <S> no no no yes yes yes no </S>*

Compute:  $P(\text{yes no no yes})$

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

Corpus:

①  $\langle S \rangle \underbrace{\langle S \rangle}_{1 \text{ time}} \text{ yes no no no no yes } \langle /S \rangle$   
 ②  $\langle S \rangle \underbrace{\langle S \rangle}_{1 \text{ time}} \text{ no no no yes yes yes no } \langle /S \rangle$

Compute:  $P(\text{yes no no yes})$

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

$P(\text{yes}   \underbrace{\langle S \rangle \langle S \rangle}_{\text{occurs two times}})$	1/2	$\langle S \rangle \langle S \rangle \text{ yes}$ $\langle S \rangle \langle S \rangle \text{ no}$
--	-----	---

Corpus:

<S> <S> yes no no no no yes </S>  
<S> <S> no no no yes yes yes no </S>

Compute:  $P(\text{yes } \text{no } \text{no } \text{yes})$

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

$P(\text{yes}   \text{<s>} \text{<s>})$	1/2	<S> <S> yes <S> <S> no
$P(\text{no}   \text{<s>} \text{yes})$	1/1	<S> yes no

## Corpus:

$\langle s \rangle \langle s \rangle \text{ yes no no no no yes } \langle /s \rangle$   
 $\langle s \rangle \langle s \rangle \text{ no no no yes yes yes no } \langle /s \rangle$

Compute:  $P(\text{yes no no yes})$

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

$P(\text{yes}   \langle s \rangle \langle s \rangle)$	1/2	$\langle s \rangle \langle s \rangle \text{ yes}$ $\langle s \rangle \langle s \rangle \text{ no}$
$P(\text{no}   \langle s \rangle \text{ yes})$	1/1	$\langle s \rangle \text{ yes no}$
$P(\text{no}   \text{ yes no})$	1/2	$\text{yes no no}$ $\text{yes no } \langle /s \rangle$

## Corpus:

$\langle s \rangle \langle s \rangle yes no no \underline{no} \underline{no} yes \langle /s \rangle$   
 $\langle s \rangle \langle s \rangle no \underline{no} \underline{no} yes yes yes no \langle /s \rangle$

Compute:  $P(yes \text{ } no \text{ } no \text{ } yes)$

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

$P(yes   \langle s \rangle \langle s \rangle)$	1/2	$\langle s \rangle \langle s \rangle yes$ $\langle s \rangle \langle s \rangle no$
$P(no   \langle s \rangle yes)$	1/1	$\langle s \rangle yes no$
$P(no   yes no)$	1/2	$yes no no$ $yes no \langle /s \rangle$
$P(yes   no no)$	? $\cancel{2/5}$	

[PollEv.com/jeyhanlau569](https://PollEv.com/jeyhanlau569)



## Corpus:

*<S> <S> yes no no no no yes </S>  
 <S> <S> no no no yes yes yes no </S>*

Compute:  $P(\text{yes no no yes})$

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

$P(\text{yes}   <\text{S}> <\text{S}>)$	1/2	<i>&lt;S&gt; &lt;S&gt; yes      &lt;S&gt; &lt;S&gt; no</i>
$P(\text{no}   <\text{S}> \text{ yes})$	1/1	<i>&lt;S&gt; yes no</i>
$P(\text{no}   \text{ yes no})$	1/2	<i>yes no no      yes no &lt;/S&gt;</i>
$P(\text{yes}   \text{ no no})$	2/5	<i>no no no      no no no      no no yes      no no no      no no yes</i>

## Corpus:

$\langle s \rangle \langle s \rangle yes no no no no \underline{no yes} \langle /s \rangle$   
 $\langle s \rangle \langle s \rangle no no no yes yes yes no \langle /s \rangle$

Compute:  $P(yes \text{ } no \text{ } no \text{ } yes)$

$$P(w_i | w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

$P(yes   \langle s \rangle \langle s \rangle)$	1/2	$\langle s \rangle \langle s \rangle yes$ $\langle s \rangle \langle s \rangle no$
$P(no   \langle s \rangle yes)$	1/1	$\langle s \rangle yes no$
$P(no   yes no)$	1/2	$yes no no$ $yes no \langle /s \rangle$
$P(yes   no no)$	2/5	no no no no no no no no yes no no no no no yes
$P(\langle /s \rangle   no yes)$	1/2	no yes $\langle /s \rangle$ no yes yes

## Corpus:

$\langle s \rangle \langle s \rangle yes no no no no yes \langle /s \rangle$   
 $\langle s \rangle \langle s \rangle no no no yes yes yes no \langle /s \rangle$

Compute:  $P(yes \text{ } no \text{ } no \text{ } yes) = \frac{1}{2} \times 1 \times \frac{1}{2} \times \frac{2}{5} \times \frac{1}{2} = 0.05$

$P(yes   \langle s \rangle \langle s \rangle)$	1/2	$\langle s \rangle \langle s \rangle yes$ $\langle s \rangle \langle s \rangle no$
$P(no   \langle s \rangle yes)$	1/1	$\langle s \rangle yes no$
$P(no   yes no)$	1/2	$yes no no$ $yes no \langle /s \rangle$
$P(yes   no no)$	2/5	no no no no no no no no yes no no no no no yes
$P(\langle /s \rangle   no yes)$	1/2	no yes $\langle /s \rangle$ yes no $\langle /s \rangle$

# Several Problems

- Language has long distance effects – need large  $n$ 
  - *The lecture/s that took place last week was/were on preprocessing.*
- Resulting probabilities are often very small
  - Use log probability to avoid numerical underflow
- What about unseen  $n$ -grams?
  - $P(w_1, w_2, \dots, w_m) = P(w_1 | < s >) \times P(w_2 | w_1) \dots$
  - Need to smooth the LM!

whole term = 0 if  $P(w_2|w_1) = 0$



# Smoothing

# Smoothing

- Basic idea: give events you've never seen before some probability  
*"pseudocounts"*
- Must be the case that  $P(\text{everything}) = 1$
- Many different kinds of smoothing
  - ▶ Laplacian (add-one) smoothing
  - ▶ Add- $k$  smoothing
  - ▶ Absolute discounting
  - ▶ Kneser-Ney
  - ▶ And others...

# ① Laplacian (Add-one) Smoothing

- Simple idea: pretend we've seen each  $n$ -gram once more than we did.

For unigram models ( $V$ = the vocabulary),  
*added one pseudo count for each word you've seen in  
size of the corpus.*

$$P_{add1}(w_i) = \frac{C(w_i) + 1}{M + |V|}$$

For bigram models,

$$P_{add1}(w_i | w_{i-1}) = \frac{C(w_{i-1} w_i) + 1}{C(w_{i-1}) + |V|}$$

# Add-one Example

*is only ever used as a prior context.  $\Rightarrow$  never add pseudo count to <s>*

**<s> the rat ate the cheese </s>**

*$\hookrightarrow$  add one coz we have to estimate it sometimes*

What's the bigram probability  $P(\text{ate} \mid \text{rat})$  under add-one smoothing?

$$= \frac{C(\text{rat ate}) + 1}{C(\text{rat}) + |V|} = \frac{2}{6}$$

$V = \{ \text{the}, \text{rat}, \text{ate}, \text{cheese}, \text{</s>} \}$

What's the bigram probability  $P(\text{ate} \mid \text{cheese})$  under add-one smoothing?

$$= \frac{C(\text{cheese ate}) + 1}{C(\text{cheese}) + |V|} = \frac{1}{6}$$

*<s> is not part of vocabulary because we never need to infer its conditional probability (e.g.  $P(<\text{s}> \mid \dots)$ ) predict it.*

Recall:  $P(\text{yes no no yes}) =$   
 $P(\text{yes} \mid \text{<s> } \text{<s>}) \times P(\text{no} \mid \text{<s> yes}) \times$   
 $P(\text{no} \mid \text{yes no}) \times P(\text{yes} \mid \text{no no}) \times$   
 $P(\text{</s>} \mid \text{no yes})$

(2)

## Add- $k$ Smoothing

- Adding one is often too much
- Instead, add a fraction  $k$
- AKA Lidstone Smoothing, or Add- $\alpha$  Smoothing

$$P_{addk}(w_i | w_{i-1}, w_{i-2}) = \frac{C(w_{i-2}, w_{i-1}, w_i) + k}{C(w_{i-2}, w_{i-1}) + k |V|}$$

- Have to choose  $k$  (Trick)

maintain the prob.  
dist.

real counts ↑

effective counts taken ↑

Prior Context = *alleged*

- 5 observed bi-grams
- 2 unobserved bi-grams

a bit smaller

<u>Alleged + Fay 73</u>		counts	unsmoothed probability	Lidstone smoothing, $\alpha = 0.1$	
				effective counts <small>pseudo counts given</small>	smoothed probability
<i>impropriety</i>	8	0.4	7.826	0.391	
<i>offense</i>	5	0.25	4.928	0.246	
<i>damage</i>	4	0.2	3.961	0.198	
<i>deficiencies</i>	2	0.1	2.029	0.101	
<i>outbreak</i>	1	0.05	1.063	0.053	
<i>infirmitiy</i>	0	0	0.097	0.005	
<i>alleged</i>	0	0	0.097	0.005	

20

1.0

20

1.0

$$0.391 \times 20$$

$$(0 + 0.1) / (20 + 7 \times 0.1)$$

$$(8 + 0.1) / (20 + 7 \times 0.1)$$

|V|

(3)

# Absolute Discounting

*take same effective counts*

- ‘Borrows’ a **fixed probability mass** from observed n-gram counts
- Redistributions it to unseen n-grams

# Absolute Discounting

- Context = *alleged*
- 5 observed bi-grams
  - 2 unobserved bi-grams

	counts	unsmoothed probability	Lidstone smoothing, $\alpha = 0.1$		Discounting, $d = 0.1$	
			effective counts	smoothed probability	effective counts (1)	smoothed probability (2)
<i>impropriety</i>	8	0.4	7.826	0.391	7.9	0.395
<i>offense</i>	5	0.25	4.928	0.246	4.9	0.245
<i>damage</i>	4	0.2	3.961	0.198	3.9	0.195
<i>deficiencies</i>	2	0.1	2.029	0.101	1.9	0.095
<i>outbreak</i>	1	0.05	1.063	0.053	0.9	0.045
<i>infirmity</i>	0	0	0.097	0.005	0.25	0.013
<i>alleged</i>	0	0	0.097	0.005	0.25	0.013

*factor*

20      1.0      20      1.0      20      1.0

8 - 0.1      (0.1 x 5) / 2      0.25 / 20

two unseen ones  
0.1 have been taken  
5 times

effective counts      total tokens

0.1      0.25      0.25

27

④

# Backoff

freq.↑ weight↑.

- Absolute discounting redistributes the probability mass **equally** for all unseen n-grams

↳ Not realistic

- Katz Backoff: redistributes the mass **based on a lower order model** (e.g. unigram)

$$P_{\text{katz}}(w_i | w_{i-1}) = \begin{cases} \frac{C(w_{i-1}, w_i) - D}{C(w_{i-1})}, & \text{for seen bigrams} \\ \alpha(w_{i-1}) \times \frac{P(w_i)}{\sum_{w_j: C(w_{i-1}, w_j) = 0} P(w_j)}, & \text{otherwise} \end{cases}$$

*What you have taken from the seen bigrams*

*sum unigram probabilities for all words that do not co-occur with context  $w_{i-1}$*   
*e.g.  $P(\text{infirmity}) + P(\text{alleged})$*

*unigram probability for  $w_i$   
e.g.  $P(\text{infirmity})$*

*weight each unseen bigram based on unigrams*

*decide how much mass we give to the unseen bigrams if  $C(w_{i-1}, w_i) > 0$*

*the amount of probability mass that has been discounted for context  $w_{i-1}$   
((0.1 x 5) / 20 in previous slide)*

# Issues with Katz Backoff

$$P_{katz}(w_i|w_{i-1}) = \begin{cases} \frac{C(w_{i-1}, w_i) - D}{C(w_{i-1})}, & \text{if } C(w_{i-1}, w_i) > 0 \\ \alpha(w_{i-1}) \times \frac{P(w_i)}{\sum_{w_j: C(w_{i-1}, w_j) = 0} P(w_j)}, & \text{otherwise} \end{cases}$$

- I can't see without my reading what's most likely next word?  
*bigram*
- $C(\text{reading}, \text{glasses}) = C(\text{reading}, \text{Francisco}) = 0$   
*Unigram*
- $C(\text{Francisco}) > C(\text{glasses})$  *in a corpus (assumed)*
- Katz backoff will give higher probability to *Francisco* *But actually doesn't make sense*

(5)

## Kneser-Ney Smoothing

- Redistribute probability mass based on the **versatility** of the lower order n-gram
- AKA “continuation probability”
- What is versatility? *how often a word co-occur with other unique word.*
  - High versatility -> co-occurs with a lot of unique words, e.g. glasses
    - men's glasses, black glasses, buy glasses, etc
  - Low versatility -> co-occurs with few unique words, e.g. francisco
    - san francisco

# Kneser-Ney Smoothing

Seen bigram ↗

$$P_{KN}(w_i | w_{i-1}) = \begin{cases} \frac{C(w_{i-1}, w_i) - D}{C(w_{i-1})}, & \text{if } C(w_{i-1}, w_i) > 0 \\ \beta(w_{i-1}) P_{cont}(w_i), & \text{otherwise} \end{cases}$$

the amount of probability mass that has been discounted for context  $w_{i-1}$

$\alpha \approx 1$

看  $w_i$  有多少种不同种类的 word 有几种不同的  $w_{i-1}$  加起来 ← All the diff  $w_i$ . count the freq. of versatility

$$P_{cont}(w_i) = \frac{|\{w_{i-1} : C(w_{i-1}, w_i) > 0\}|}{\sum_{w_j} |\{w_{j-1} : C(w_{j-1}, w_j) > 0\}|}$$

- Intuitively the numerator of  $P_{cont}$  counts the number of unique  $w_{i-1}$  that co-occurs with  $w_i$
- High continuation counts for glasses
- Low continuation counts for Francisco

$C(\text{word type of } w_{i-1})$

(6)

# Interpolation

- A better way to combine different orders of  $n$ -gram models

$$\text{Sum} = \text{unigram} + \text{bigram} + \text{trigram}$$

- Interpolated trigram model:

$$P_{IN}(w_i | w_{i-1}, w_{i-2}) = \lambda_3 P_3(w_i | w_{i-2}, w_{i-1}) + \lambda_2 P_2(w_i | w_{i-1}) + \lambda_1 P_1(w_i)$$

Learned based on held out data

$$\lambda_1 + \lambda_2 + \lambda_3 = 1$$

# A Final Word

- $N$ -gram language models are a simple but effective way to capture the predictability of language
- Can be trained in an unsupervised fashion, scalable to large corpora
- Require smoothing to be effective
- Modern language models uses neural networks (lecture 8)

# Reading

- E18 Chapter 6 (skip 6.3)

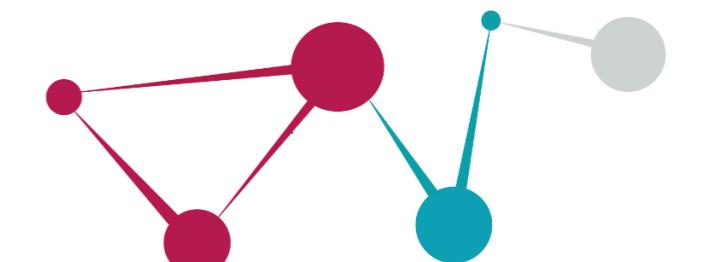
<b>Events Confirmed</b>	<b>Wee k</b>	<b>Date, Time, Location</b>
<b>Summerfest Welcome Back Expo</b>	<b>1</b>	1-Mar-22, 2-5pm 2-Mar-22, 2-5pm 3-Mar-22, 2-5pm 4-Mar-22, 2-5pm
<b>Welcome event with DSSS &amp; HackMelbourne</b>	<b>2</b>	7-Mar-22, 11:30am
<b>Google Career Discussion</b>	<b>2</b>	8-Mar-22, 5 - 6:30pm
<b>Spaghetti Bridge Building Comp Collab</b>	<b>2</b>	11-Mar-22, 3-5:30pm
<b>Robogals collab industry week</b>	<b>3</b>	<b>this whole week</b>
<b>Google Women in Tech clubs session</b>	<b>3</b>	16-Mar-22, 5-5:45pm
<b>Macquarie Grad Program Event</b>	<b>3</b>	<b>TBD</b>

## **Events planned (wk4 onwards):**

- Annual Hackathon
- Workshops for CV and job seeking
- High Tea networking
- Industry Panels
- Free lunch on campus
- Study sessions
- Tech-related clubs collaboration
- SWOT cake
- ...

## **Other opportunities:**

- Weekly job posts
- Tech-related events/competitions
- Volunteering opportunities
- ...



**WOMEN IN TECH**

The University of Melbourne

**JOIN US HERE!**



# JOIN DATA SCIENCE STUDENT SOCIETY TODAY

TOMORROW'S DATA SCIENTIST  
STARTS WITH US



**FREE MEMBERSHIP/  
CLUB MAILING LIST**

OR VISIT

<https://tinyurl.com/dscubed-unimelb-2022>



[www.facebook.com/dscubed.unimelb](http://www.facebook.com/dscubed.unimelb)



[www.instagram.com/dscubed.unimelb/](http://www.instagram.com/dscubed.unimelb/)



[www.linkedin.com/company/data-science-student-society](http://www.linkedin.com/company/data-science-student-society)