# Natural Language Processing Language modeling

Sofia Serrano sofias6@cs.washington.edu

Credit to Yulia Tsvetkov and Noah Smith for slides

# Announcements

- Quiz 2 will be released on Canvas at the end of lecture today
  - Available from 2:20pm today through 2:20pm Thursday (tomorrow)
  - 5 questions, 10 minutes
- A1 is due at 11:59pm this Friday
  - For details on how to tag and submit your assignment, see Leo's tutorial video linked in Sunday's Ed announcement
- We're holding extra office hours this week
  - See the office hour schedule <u>google doc</u> from the Sunday Ed announcement, also linked on the course website under "Announcements"

#### Language models in the news (these days, ChatGPT)

Hi, I'm writing an article about you and all of the cool things you can do. Could you say hello to my readers and give us a quick self-introduction? I'd like it to be friendly and casual, include a silly joke, and end with a note about how great the rest of my article is.



Hello readers! My name is Assistant, and I'm a large language model trained by OpenAI. I'm here to help answer your questions and provide you with information on a wide variety of topics. I'm a bit of a jokester, so here's a quick one for you: Why was the math book sad? Because it had too many problems! But seriously, I'm excited to be featured in your article and I hope you enjoy learning more about me and all the cool things I can do. And trust me, the rest of the article is going to be great!

Image taken from Springboard

ሰ ጥ

# What is language modeling?

• Language Modeling is the task of predicting what word comes next



#### We use language models every day



#### We use language models every day

# Google

what is the			Ŷ
what is the <b>weather</b> what is the <b>meanin</b> what is the <b>dark we</b> what is the <b>dark we</b> what is the <b>doomso</b> what is the <b>doomso</b> what is the <b>weather</b> what is the <b>keto die</b> what is the <b>america</b> what is the <b>speed o</b> what is the <b>bill of ri</b>	g of life bb day clock r today et an dream of light ghts		
	Google Search	I'm Feeling Lucky	

## Language modeling



# Consider this botched small talk

"I'm often misunderstood at my supermarket in Sussex, not because of my accent but because I tend to deviate from the script.

Cashier: Hello, how are you this evening?

Me: Has your house ever been burgled?

Cashier: What?

Me: Your house- has anyone ever broken into it and stolen things?"

- David Sedaris, *Calypso* 

# Language models play the role of ...

- a judge of grammaticality
  - e.g., should prefer "The boy runs." to "The boy run."
- a judge of semantic plausibility
  - e.g., should prefer "The woman spoke." to "The sandwich spoke."
- an enforcer of stylistic consistency
  - e.g., should prefer "Hello, how are you this evening? Fine, thanks, how are you?" to "Hello, how are you this evening? Has your house ever been burgled?"
- a repository of knowledge (?)
  - e.g., "Barack Obama was the 44th President of the United States"

Note that this is very difficult to guarantee!

# The language modeling problem

- Assign a probability to every sentence (or any string of words)
  - finite vocabulary (e.g. words or characters) {*the*, *a*, *telescope*, ...}

10

- infinite set of sequences
  - a telescope STOP
  - a STOP
  - the the the STOP
  - I saw a woman with a telescope STOP
  - STOP
  - ...

# The language modeling problem

- Assign a probability to every sentence (or any string of words)
  - finite vocabulary (e.g. words or characters)
  - infinite set of sequences

$$\sum_{\mathbf{e}\in\Sigma^*} p_{\mathrm{LM}}(\mathbf{e}) = 1$$
$$p_{\mathrm{LM}}(\mathbf{e}) \ge 0 \quad \forall \mathbf{e}\in\Sigma^*$$

# Formalizing our definition

• Language Modeling is the task of predicting what word comes next



• More formally: given a sequence of words  $x^{(1)}$ ,  $x^{(2)}$ , ...  $x^{(t)}$ compute the probability distribution of the next word  $x^{(t+1)}$ where  $x^{(t+1)}$  can be any word in the vocabulary V={  $w_1, w_2, ..., w_{|V|}$ }



 $\mathscr{V} = \{\text{permit, reject}\}$  Our event space is  $\mathscr{V}^*$  with <eos> at end Our r.v. is X



 $\mathscr{V} = \{\text{permit, reject}\}$  Our event space is  $\mathscr{V}^*$  with <eos> at end Our r.v. is X





 $\mathscr{V} = \{\text{permit, reject}\}$  Our event space is  $\mathscr{V}^*$  with <eos> at end Our r.v. is X



15



 $\mathscr{V} = \{\text{permit, reject}\}$  Our event space is  $\mathscr{V}^*$  with <eos> at end Our r.v. is X



16



 $\mathscr{V} = \{\text{permit, reject}\}$  Our event space is  $\mathscr{V}^*$  with <eos> at end Our r.v. is X





 $\mathcal{V} = \{\text{permit, reject}\}$  Our event space is  $\mathcal{V}^*$  with  $\langle \cos \rangle$  at end Our r.v. is X

We'll say that X is the distribution of utterances this person produces, and we want to estimate X.





 $\mathscr{V}$  = {permit, reject} Our event space is  $\mathscr{V}$ \* with <eos> at end Our r.v. is X What is p(X = reject permit <eos>)?



# Language Modeling

• If we have some text, then the probability of this text (according to the Language Model) is:

$$P(\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(T)}) = P(\boldsymbol{x}^{(1)}) \times P(\boldsymbol{x}^{(2)} | \boldsymbol{x}^{(1)}) \times \dots \times P(\boldsymbol{x}^{(T)} | \boldsymbol{x}^{(T-1)}, \dots, \boldsymbol{x}^{(1)})$$
$$= \prod_{t=1}^{T} P(\boldsymbol{x}^{(t)} | \boldsymbol{x}^{(t-1)}, \dots, \boldsymbol{x}^{(1)})$$

This is what our LM provides



 $\mathscr{V}$  = {permit, reject} Our event space is  $\mathscr{V}$  with <eos> at end Our r.v. is X What is p(X = reject permit <eos>)?





 $\mathscr{V}$ = {permit, reject} Our event space is  $\mathscr{V}$ \* with <eos> at end Our r.v. is X

What is p(X = reject permit <eos>)? Use chain rule of probability.





 $\mathscr{V}$ = {permit, reject} Our event space is  $\mathscr{V}$ \* with <eos> at end Our r.v. is X

p(X = reject permit <eos>) = p(reject | [start]) \* p(permit | [start] reject) \*





 $\mathscr{V} = \{\text{permit, reject}\}$  Our event space is  $\mathscr{V}^*$  with <eos> at end Our r.v. is X





 $\mathcal{V} = \{\text{permit, reject}\}$  Our event space is  $\mathcal{V}^*$  with <eos> at end Our r.v. is X Note: as long as p(child of node) > 0 for each node and  $\sum p(child of node) = 1$  for child each (non-eos) node, then  $\sum p(path) = 1$ path in tree [start] reject permit < eos >reject permit permit reject <eos> <eos> permit <eos> reject permit <eos> reject permit <eos> reject permit <eos> reject

 $p(how are you this evening ? has your house ever been burgled ? STOP) = <math>10^{-15}$  $p(how are you this evening ? fine , thanks , how about you ? STOP) = <math>10^{-9}$ 

# **Motivation**

• Speech recognition: we want to predict a sentence given acoustics



### **Motivation**

• Speech recognition: we want to predict a sentence given acoustics

the station signs are indeed in english the station signs are in deep in english the stations signs are in deep in english the station signs are in deep into english the station 's signs are in deep in english the station signs are in deep in the english the station 's signs are indeed in english the station signs are indians in english the station signs are indian in english the stations signs are indians in english the stations signs are indians and english

-14725 -14732 -14735

- -14739 -14740
- -14741
- -14760
- -14790
- -14799
  - -14807
  - -14815

# **Motivation**

- Machine translation
  - p(strong winds) > p(large winds)
- Spelling correction
  - The office is about fifteen minuets from my house
  - p(about fifteen minutes from) > p(about fifteen minuets from)

- Speech recognition
  - p(I saw a van) >> p(eyes awe of an)

• Summarization, question-answering, handwriting recognition, OCR, etc.

# A trivial model

- Assume we have **n** training sentences
- Let  $x_1, x_2, ..., x_n$  be a sentence, and  $c(x_1, x_2, ..., x_n)$  be the number of times it appeared in the training data.
- Define a language model:

$$p(x_1,\ldots,x_n) = \frac{c(x_1,\ldots,x_n)}{N}$$

• No generalization!

#### Sentence/paragraph/book probability

$$P(\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(T)}) = P(\boldsymbol{x}^{(1)}) \times P(\boldsymbol{x}^{(2)} | \boldsymbol{x}^{(1)}) \times \dots \times P(\boldsymbol{x}^{(T)} | \boldsymbol{x}^{(T-1)}, \dots, \boldsymbol{x}^{(1)})$$
$$= \prod_{t=1}^{T} P(\boldsymbol{x}^{(t)} | \boldsymbol{x}^{(t-1)}, \dots, \boldsymbol{x}^{(1)})$$

P(its water is so transparent that the) =

P(its)	×
P(water   its)	×
P(is   its water)	×
P(so   its water is)	×
P(transparent   its water is so)	×
	×

P(the | its water is so transparent that)  $\rightarrow$  How to estimate?

#### 32

# **Markov assumption**

- We make the Markov assumption:  $\mathbf{x}^{(t+1)}$  depends only on the preceding n-1 words
  - Markov chain is a "...stochastic model describing a sequence of possible events in which the probability of each event depends only on the state attained in the previous event."

$$P(\boldsymbol{x}^{(t+1)} | \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(1)}) = P(\boldsymbol{x}^{(t+1)} | \boldsymbol{x}^{(t)}, \dots, \boldsymbol{x}^{(t-n+2)})$$
  
n-1 words

assumption



Andrei Markov

# **Markov assumption**

#### P(the | its water is so transparent that) $\equiv$ P(the | transparent that)

Andrei Markov

or maybe even

P(the | its water is so transparent that)  $\equiv$  P(the | that)

33

"I have a dog whose name is Lucy. I have two cats, they like playing with Lucy."

• Question: How to learn a Language Model?

- Question: How to learn a Language Model?
- Answer (pre- Deep Learning): learn an *n-gram* Language Model!

"I have a dog whose name is Lucy. I have two cats, they like playing with Lucy."

• Definition: An n-gram is a chunk of n consecutive words.

- Definition: An n-gram is a chunk of n consecutive words.
  - unigrams: {I, have, a, dog, whose, name, is, Lucy, two, cats, they, like, playing, with}

- Definition: An n-gram is a chunk of n consecutive words.
  - unigrams: {I, have, a, dog, whose, name, is, Lucy, two, cats, they, like, playing, with}
  - bigrams: {I have, have a, a dog, dog whose, ..., with Lucy}

- Definition: An n-gram is a chunk of n consecutive words.
  - unigrams: {I, have, a, dog, whose, name, is, Lucy, two, cats, they, like, playing, with}
  - bigrams: {I have, have a, a dog, dog whose, ..., with Lucy} have cats

- Definition: An n-gram is a chunk of n consecutive words.
  - unigrams: {I, have, a, dog, whose, name, is, Lucy, two, cats, they, like, playing, with}
  - bigrams: {I have, have a, a dog, dog whose, ..., with Lucy} have cats 🗙

- Definition: An n-gram is a chunk of n consecutive words.
  - unigrams: {I, have, a, dog, whose, name, is, Lucy, two, cats, they, like, playing, with}
  - bigrams: {I have, have a, a dog, dog whose, ..., with Lucy}
  - trigrams: {I have a, have a dog, a dog whose, ..., playing with Lucy}

- Definition: An n-gram is a chunk of n consecutive words.
  - unigrams: {I, have, a, dog, whose, name, is, Lucy, two, cats, they, like, playing, with}
  - bigrams: {I have, have a, a dog, dog whose, ..., with Lucy}
  - trigrams: {I have a, have a dog, a dog whose, ..., playing with Lucy}
  - four-grams: {I have a dog, ..., like playing with Lucy}
  - 0 ...

- $w_1 a$  unigram
- $w_1 w_2 a bigram$
- $w_1 w_2 w_3$  a trigram
- $w_1 w_2 \dots w_n$  an n-gram

"I have a dog whose name is Lucy. I have two cats, they like playing with Lucy."

- Question: How to learn a Language Model?
- Answer (pre- Deep Learning): learn an *n-gram* Language Model!

• Idea: Collect statistics about how frequent different n-grams are and use these to predict next word

# unigram probability

"I have a dog whose name is Lucy. I have two cats, they like playing with Lucy."

45

- corpus size m = 17
- P(Lucy) = 2/17; P(cats) = 1/17

• Unigram probability: 
$$P(w) = \frac{count(w)}{m} = \frac{C(w)}{m}$$

# bigram probability

"I have a dog whose name is Lucy. I have two cats, they like playing with Lucy."

0

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

$$P(have | I) = \frac{P(I have)}{P(I)} = \frac{2}{2} = 1$$

$$P(two | have) = \frac{P(have two)}{P(have)} = \frac{1}{2} = 0.5$$

$$P(eating | have) = \frac{P(have eating)}{P(have)} = \frac{0}{2} = 0$$

$$P(w_2|w_1) = \frac{C(w_1, w_2)}{\sum_{w} C(w_1, w)} = \frac{C(w_1, w_2)}{C(w_1)}$$

# trigram probability

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

$$P(a \mid I \text{ have}) = \frac{C(I \text{ have } a)}{C(I \text{ have})} = \frac{1}{2} = 0.5$$

$$P(w_3 \mid w_1 w_2) = \frac{C(w_1, w_2, w_3)}{\sum_w C(w_1, w_2, w)} = \frac{C(w_1, w_2, w_3)}{C(w_1, w_2)}$$

$$P(\text{several} \mid I \text{ have}) = \frac{C(I \text{ have several})}{C(I \text{ have})} = \frac{0}{2} = 0$$

#### n-gram probability

"I have a dog whose name is Lucy. I have two cats, they like playing with Lucy."

48

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

$$P(w_i | w_1, w_2, ..., w_{i-1}) = \frac{C(w_1, w_2, ..., w_{i-1}, w_i)}{C(w_1, w_2, ..., w_{i-1})}$$

#### **First-order Markov process**

#### Chain rule

$$p(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) =$$
$$p(X_1 = x_1) \prod_{i=2}^n p(X_i = x_i \mid X_1 = x_1, \dots, X_{i-1} = x_{i-1})$$

#### **First-order Markov process**

#### Chain rule

$$p(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) =$$
$$p(X_1 = x_1) \prod_{i=2}^n p(X_i = x_i \mid X_1 = x_1, \dots, X_{i-1} = x_{i-1})$$

50

Markov assumption

$$= P(X_1 = x_1) \prod_{i=2}^{n} P(X_i = x_i | X_{i-1} = x_{i-1})$$

#### Second-order Markov process:

• Relax independence assumption:

$$p(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) =$$

$$p(X_1 = x_1) \times p(X_2 = x_2 \mid X_1 = x_1)$$

$$\times \prod_{i=3}^n p(X_i = x_i \mid X_{i-2} = x_{i-2}, X_{i-1} = x_{i-1})$$

#### Second-order Markov process:

• Relax independence assumption:

$$p(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) =$$

$$p(X_1 = x_1) \times p(X_2 = x_2 \mid X_1 = x_1)$$

$$\times \prod_{i=3}^n p(X_i = x_i \mid X_{i-2} = x_{i-2}, X_{i-1} = x_{i-1})$$

• Simplify notation:

$$x_0 = *, x_{-1} = *$$

# 3-gram LMs

- A trigram language model contains
  - a vocabulary V
  - a non negative parameters q(w|u,v) for every trigram, such that

$$w \in \mathcal{V} \cup \{\text{STOP}\}, \ u, v \in \mathcal{V} \cup \{*\}$$

• the probability of a sentence  $x_1, ..., x_n$ , where  $x_n = STOP$  is

$$p(x_1, \dots, x_n) = \prod_{i=1}^n q(x_i \mid x_{i-1}, x_{i-2})$$

#### Example

p(the dog barks STOP) =

#### Example

#### $p(\text{the dog barks STOP}) = q(the \mid *, *) \times$

#### Example

 $p(\text{the dog barks STOP}) = q(the \mid *, *) \times$   $q(dog \mid *, the) \times$   $q(barks \mid the, dog) \times$   $q(STOP \mid dog, barks) \times$ 

56

# Berkeley restaurant project sentences

- can you tell me about any good cantonese restaurants close by
- mid priced that food is what i'm looking for
- tell me about chez pansies
- can you give me a listing of the kinds of food that are available
- i'm looking for a good place to eat breakfast
- when is caffe venezia open during the day

#### Raw bigram counts (~1000 sentences)

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

58

#### **Bigram probabilities**

$$P(w_{i} | w_{i-1}) = \frac{C(w_{i-1}, w_{i})}{C(w_{i-1})}$$
$$P(w_{1}, w_{2}, ..., w_{n}) \equiv \prod_{i} P(w_{i} | w_{i-1})$$

i	want	to		eat		chinese		fe	boc		lunch	spend
2533	927	2417	7	746		158		1	.093		341	278
	:	mont	ta			+	ahina		food		hunch	amand
	1	want	10		ea	u	cnine	ese	1000		lunch	spend
i	0.002	0.33	0		0.	0036	0		0		0	0.00079
want	0.0022	0	0.6	0.66		0011	0.006	55	0.000	55	0.0054	0.0011
to	0.00083	0	0.0	0.0017		28	0.000	)83	0		0.0025	0.087
eat	0	0	0.0	0.0027			0.021	1	0.002	27	0.056	0
chinese	0.0063	0	0		0		0		0.52		0.0063	0
food	0.014	0	0.0	)14	0		0.000	)92	0.003	37	0	0
lunch	0.0059	0	0		0		0		0.002	29	0	0
spend	0.0036	0	0.0	0036	0		0		0		0	0

59

# **Bigram estimates of sentence probability**

60

- P(<s> i want chinese food </s>) =
  P(i|<s>)
- x P(want|i)
- x P(chinese|want)
- x P(food|chinese)
- x P(</s>|food)

. . .

$$P(w_{i} | w_{i-1}) = \frac{C(w_{i-1}, w_{i})}{C(w_{i-1})}$$
$$P(w_{1}, w_{2}, \dots, w_{n}) \equiv \prod_{i} P(w_{i} | w_{i-1})$$

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

# What can we learn from bigram estimates?

61

P(to|want) = 0.66

P(chinese want	= 0.0065
P(eat to)	= 0.28
P (i  <s>)</s>	= 0.25
P(food to)	= 0.0
P(want spend)	= 0.0

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0



Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives

] gram

2 gram Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her

- 1 Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives
- 2 gram

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions

- -To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have
- gram -Hill he late speaks; or! a more to leg less first you enter
  - –Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.
- gram –What means, sir. I confess she? then all sorts, he is trim, captain.
- 3 gram

gram

- –Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.
- m –This shall forbid it should be branded, if renown made it empty.

-King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;

-It cannot be but so.