

CS 371 N Lecture 9

Language Modeling

Announcements

- A1 back
- A2 due Thurs
- Bias due Thurs
- A3 out Thurs
- Heng Ji talk Friday 11am 6302

Recap

Skip-gram $P(\text{context}=y | \text{word}=x) =$

$$\frac{e^{\bar{v}_x \cdot \bar{c}_y}}{\sum_{y'} e^{\bar{v}_x \cdot \bar{c}_{y'}}$$

Bias, multilingual, other methods

Today - Language Modeling

- N-gram LMs

- Neural LMs

- RNN LMs (brief!)

LM "Autocomplete" / predictive text

Predict the next word w_i given words that came before w_1, \dots, w_{i-1}

$$P(w_i | w_1 \dots w_{i-1})$$

n words

Prob of a sequence \bar{w} :

$$P(\bar{w}) = \prod_{i=1}^n P(w_i | w_1 \dots w_{i-1})$$

"What is the prob. of seeing this utterance?"

- Fill-in-the-blank (cloze)
- Maintain uncertainty over upcoming words
- Grammatical error correction:

\bar{w}, \bar{w}' one edit applied

$P(\bar{w}') \stackrel{?}{>} P(\bar{w}) \Rightarrow$ more grammatical?

N-gram Language Modeling

$$P(\bar{w}) = P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2) \\ \uparrow P(w_4 | w_1, w_2, w_3) \dots$$

LM in general

Suppose we have a long book

$$P(w_{100000} | w_1 \dots w_{99999})$$

n -gram LM: only look at past
 $n-1$ words

$$n=3: \quad P(\bar{w}) = \prod_{i=1}^n P(w_i | w_{i-n+1} \dots w_{i-1})$$

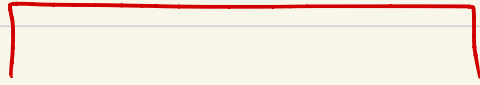
$$P(\bar{w}) = P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2) \\ P(w_4 | w_2, w_3) \dots \\ \prod_{i=1}^n w_i!$$

$n=2:$

$$P(\bar{w}) = P(w_1) P(w_2 | w_1) P(w_3 | w_2) \dots$$

Ex

I saw the dog



$n=5$

run swim
walk--
food

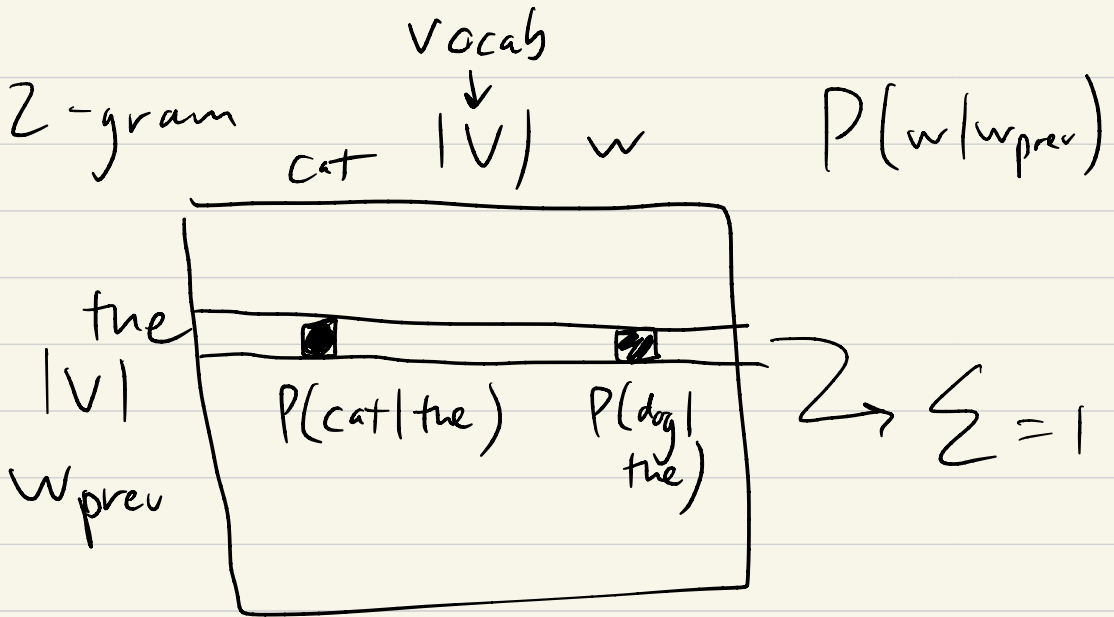
$n=2$: hates

My dog hates cats

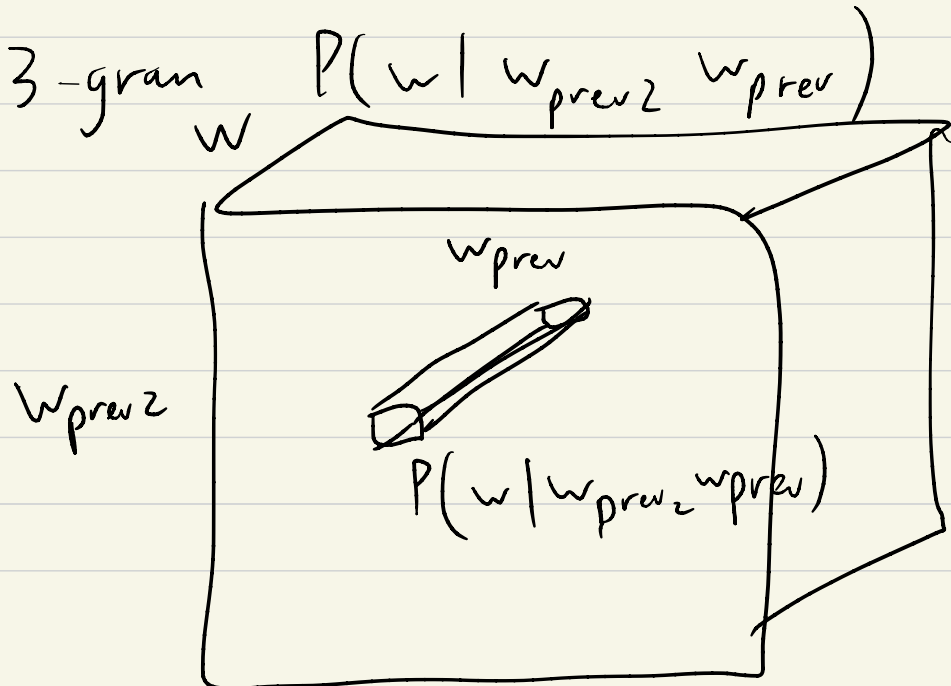
n-gram LM parameterization

Explicitly store probs for words
based on counts in the data

Categorical distribution



one row per word



These probs are sparse
actual size $\ll |V|^3$

$n \sim 7$ is doable

Parameter estimation

Count + normalize over a big corpus

the cat		
the cat	$P(\text{cat} \text{the}) =$	$\frac{\text{count}(\text{the cat})}{\text{count}(\text{the})} = \frac{2}{9}$
the dog	$P(\text{dog} \text{the}) =$	$\frac{1}{9}$
the snake	snake	$\frac{1}{9}$

Maximum-Likelihood Estimation

These probs maximize LL of
this data

Smoothing "I'd hate to go to Maui"

Suppose we have $n=5$

$$P(\text{Maui} | \text{hate to go to}) = \frac{\text{Count}(\text{hate to go to Maui})}{\text{Count}(\text{hate to go to})}$$

Solutions

- Give every 5 gram count 1

"hate to go to" seen 5 times

+ 101 "fake" counts

- Getting rid of sparsity

Smoothing a 4-gram model

$$P_4(w_i | w_{i-3} w_{i-2} w_{i-1})$$

normalized
raw
counts

$$= \lambda P_4^{\text{raw}}(w_i | w_{i-3} w_{i-2} w_{i-1})$$

$$+ (1-\lambda) P_3(w_i | w_{i-2} w_{i-1})$$

$$\rightarrow P_3(w_i | w_{i-2} w_{i-1}) =$$
$$\lambda P_3^{\text{raw}}(w_i | w_{i-2} w_{i-1})$$
$$+ (1-\lambda) P_2 \dots$$

$P_i(w) > 0$ for all $w \in V$

$\lambda = \text{hyperparameter} \in [0, 1]$

Key ideas

We can learn from a lot of data

We can model a distribution over next words given (unlimited) context

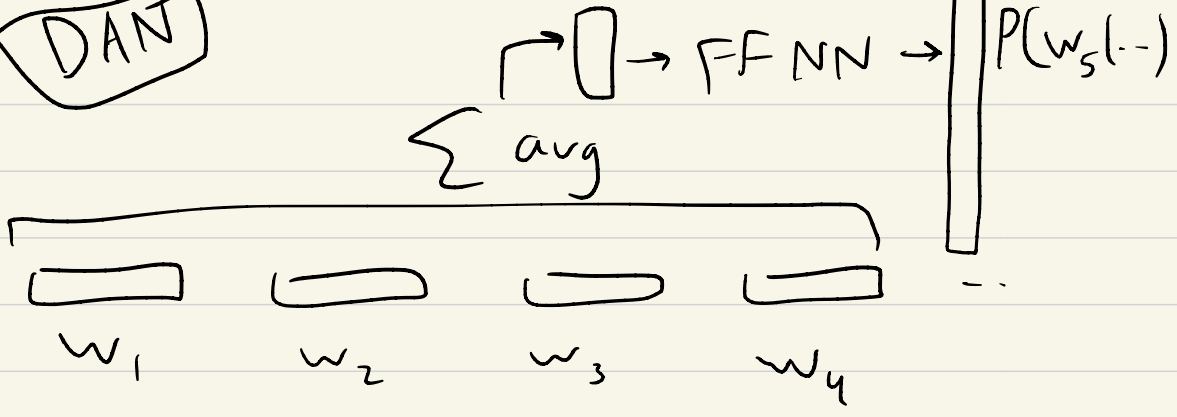
Neural Language Modeling

$P(\underline{w}_i | w_1 \dots w_{i-1}) \rightarrow$ model w/a neural net

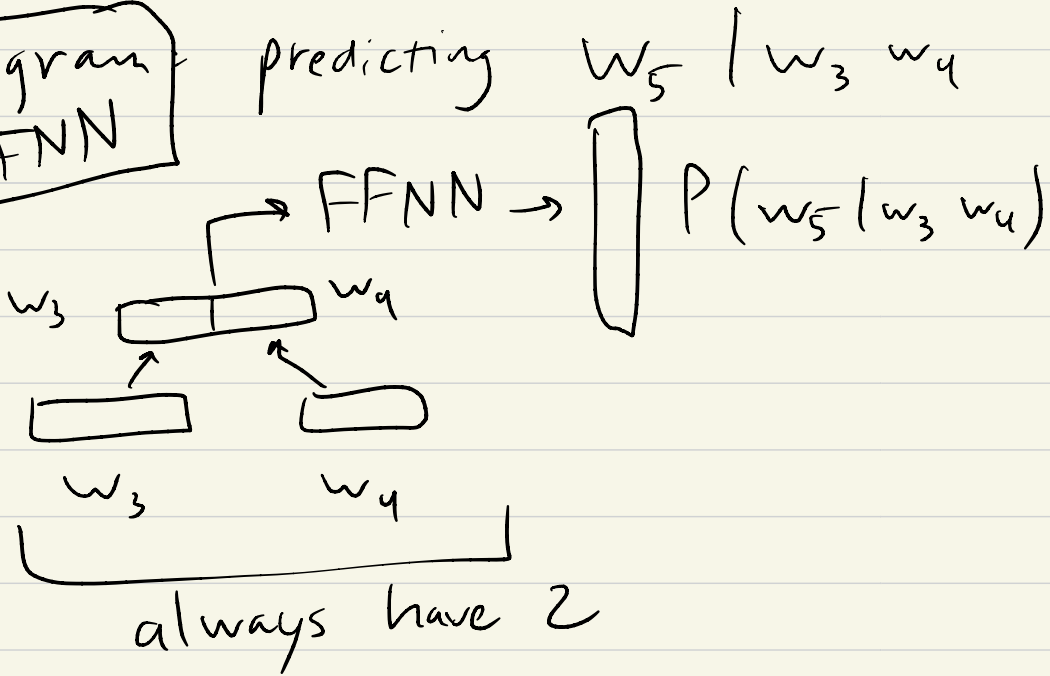
Two choices: (1) Make the n-gram approximation; (2) don't

Predict w_i based on $w_1 \dots w_{i-1}$
"sentence"

DAN



**n-gram
FFNN**



DAN Advantages:

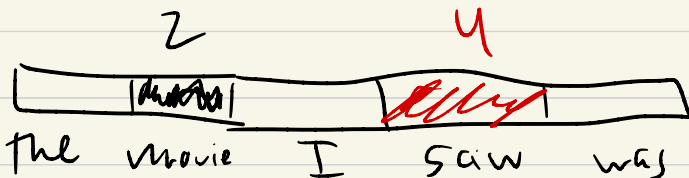
- ① Fewer params ($n=2$)
- ② Look at all the context

Disadvantages

- ① No notion of order

FFNN ① Keeps order (over a shorter context)
↳ negation!

Parameters: it's inefficient to completely model order





I saw the movie which was
FFNN distinguishes these

Solutions

- ① RNNs (recurrent neural nets):
process sequences "uniformly"
- ② Transformers: DAN + ordering
approach

RNN: maintains a "hidden state"

\bar{h} after processing n words

$$\bar{h}_n = NN(w_1, \dots, w_n)$$

$$P(w_{n+1} | w_1 \dots w_n) = \text{softmax}(U \bar{h}_n)$$

U : $|V| \times d$
matrix

Key RNN idea:

$$\bar{h}_{n+1} = \underline{\text{NN}}(\bar{h}_n, w_{n+1})$$

easy to update