

# CS371 N Lecture 9

## Language Modeling

### Announcements

- A1 + response back today
- A2 due Thurs
- Bias response due Thurs
- A3 out Thurs

### Recap Binary classification

→ multiclass

→ neural multiclass

How do we tackle "real" problems?

↳ some of these involve generation

↳ some problems are just much harder classification problems

- Today
- Language modeling
  - N-gram LMs
  - Neural LMs
  - Next time: Transformer

Language Modeling      Model distribution  $P(\bar{w})$

"Autocomplete" / predictive text

Predict the next word  $w_i$  after  
a prefix  $w_1 \dots w_{i-1}$

$$P(w_i | w_1 \dots w_{i-1}) \quad w_i \in \text{vocab } V$$

Prob of a sentence of  $n$  words:

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

What can these do?

- Anomaly detection: recognize something "weird"
- Train on one author & recognize their style
- Grammatical error correction: wrong sentence  $\bar{w}'$  has lower prob than  $\bar{w}$

## N-gram Language Modeling

In general:  $P(\bar{w}) = P(w_1) P(w_2 | w_1)$   
 $P(w_3 | w_1 w_2) \dots$   
 $P(w_{10} | w_1 w_2 \dots w_9)$

n-gram LM: only look at previous n-1 words

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

Bigram (2-gram):

$$P(\bar{w}) = P(w_1) P(w_2 | w_1) P(w_3 | w_2) \underbrace{P(w_4 | w_3) \dots}_{\downarrow}$$

probability of next  
word given prev.  
word

Count-based  
n-gram LMs

Like skip-gram!

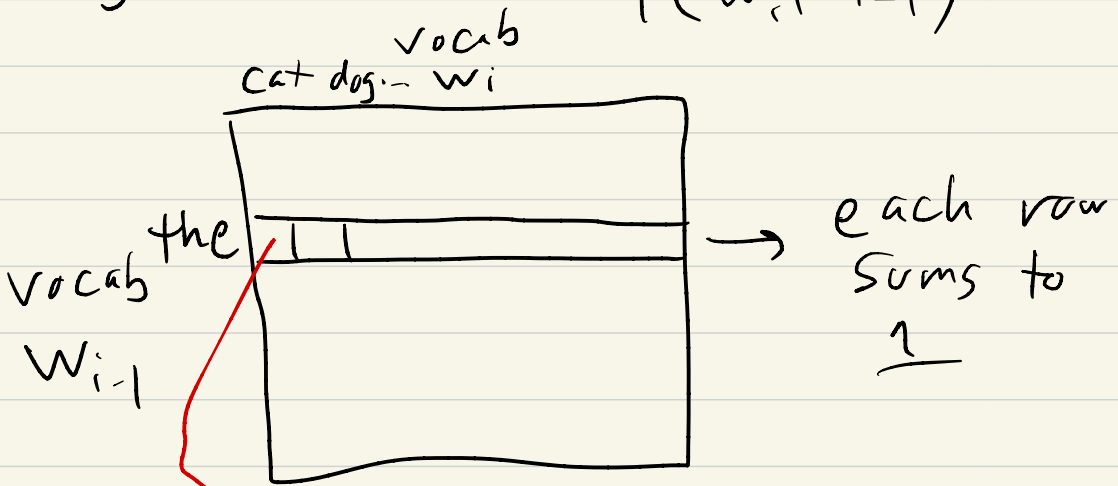
No neural nets

Explicitly represent each n-gram  
probability with a categorical dist.  
estimated from data

Z-gram LM:

next | curr

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



→  $P(\text{cat} | \text{the})$ : how likely to have cat after the?

Ex I saw the dog —

$n=2$  (bigram):  $P(w_i | \text{dog})$

What comes next? Verb (ran), noun (treat), ...

Bigram captures these okay?

⇒ ran X (wrong conjugation)



# Parameter estimation

Count + normalize over a large corpus  
(bigram)

the cat  
the cat  
the dog  
the snake

corpus

$$P(\text{cat} | \text{the}) = \frac{2}{4}$$

$$P(\text{dog} | \text{the}) = \frac{1}{4}$$

$$P(\text{snake} | \text{the}) = \frac{1}{4}$$

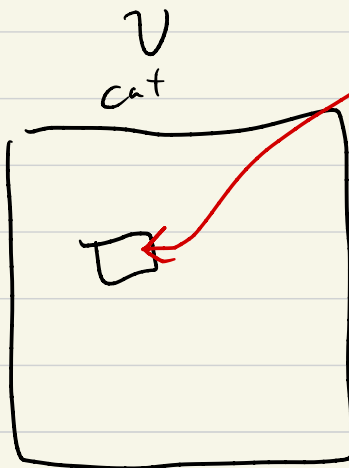
$$P(\text{the} | \text{the}) = 0$$

These parameters maximize dataset  
log likelihood

Our "model"

$V$

the



matrix instead  
of  $\bar{w}$

# Smoothing

For  $n \geq 3$ , lots of 0s in params

$P(\text{Mauil love going to}) = 0?$

If not observed

Implementing n-gram models in reality involves dealing with this

Backoff, discounting  $\swarrow$  trigram

Suppose we have  $P_3(w_i | w_{i-2} w_{i-1})$  estimated from data with MLE

$$\begin{aligned} P_{3, \text{discounted}}(w_i | w_{i-2} w_{i-1}) \\ &= \lambda P_3(w_i | w_{i-2} w_{i-1}) \\ &+ (1-\lambda) P_{2, \text{discounted}}(w_i | w_{i-1}) \end{aligned}$$



Interpolate between bigram + trigram

$$P_{2, \text{discounted}} = \lambda P_2 + (1 - \lambda) P_1$$

$P_1$  has  $> 0$  prob for every word

$\Rightarrow$  no more zeroes

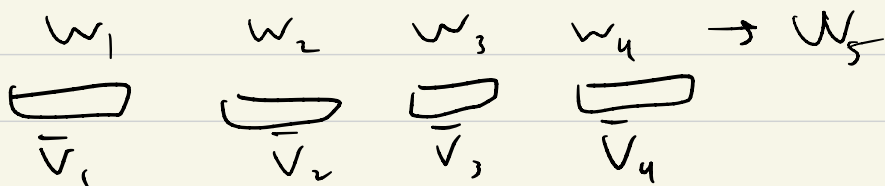
## Neural Language Modeling

$P(w_i | w_1, \dots, w_{i-1}) \Rightarrow$  model w/ a neural net

Choices: ① Do we do n-gram modeling and restrict to n-1 context words? OR look at whole sequence?

② What architecture?

Represent each token with vector



DAN  $\vec{v} = \frac{1}{n} \sum_{i=1}^n \vec{v}_i$

$$P(w_i) = \text{FFNN}(\vec{v})$$

FFNN  $\vec{v} = \text{concat}(\underbrace{v_{i-n+1} \dots v_{i-1}}_{n-1 \text{ words}})$

$n$ -gram  
↓

(position-sensitive: model the last token vs. second-to-last differently)

# DAN

# FFNN

Advantages

- Dim of  $\mathbf{V}$  is fixed

- Considers  $n-1$  words in order

- Considers whole context

Disadvantages

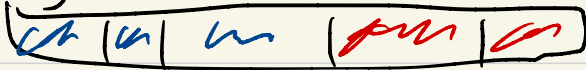
- Ignores order

- No translational invariance



The dog is wagging its —

The dog is happily wagging its —



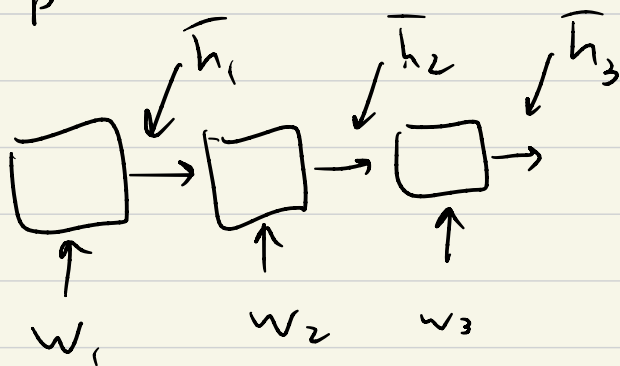
↓  
multiply by  $\mathbf{V}$   $5d_1 \times d_2$  matrix

Solutions (1) Transformer (next time)

(2) RNN - recurrent neural network

RNN forms a state vector from processing  $n$  words

Updates it with the  $n+1$  word



All  $\square$  here shared parameters

Problem  $O(n)$  sequential computation

"forget" early things in the sequence