

CS 378 Lecture 15: Language Modeling, RNNs

Today

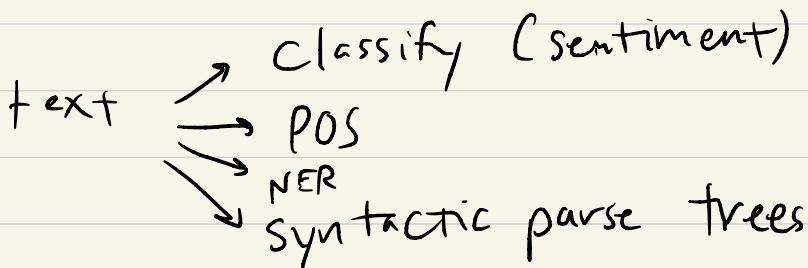
- Intro to language modeling
- N-gram LMs
- Neural LMs
- (Start) RNN LMs (A4)

Announcements

- Midterm, A3 grading
- A4
- FP custom proposals (optional!)

due Monday
(changed)

Recap (so far)



text → label or structure

Next few weeks: text \rightarrow text

(example: machine translation) seq2seq
models

dialogue, summarization

... everything?

Today: Language modeling

"autocomplete", predictive text

predict the next word given words
that came before it

Technique: recurrent neural networks
(RNNs) (+ Transformers)

sequence models

Language Modeling

Distribution $P(\bar{w})$ over (grammatical, well-formed, natural) sentences in a language

\bar{w} is a sequence of words

Why LM?

Grammatical error correction.

You give me \bar{w}

I fix some errors by finding
 \bar{w}' s.t. $P(\bar{w}') \geq P(\bar{w})$

Machine translation:

Sentence \bar{w}_S in source language

$\rightarrow \bar{w}_{t,1}$ two candidates -

Check if $P(\bar{w}_{t,1}) \geq P(\bar{w}_{t,2})$
return higher

N-gram language modeling

By the chain rule of probability:

$$P(\bar{w}) = P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2) \\ P(w_4 | w_1, w_2, w_3) \dots \\ P(w_n | w_1, \dots, w_{n-1})$$

(not Markov assumption)

If we make an assumption:

only depend on past $n-1$ words

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

2-gram model: the cat ran

$$P_2(\bar{w}) = P(\text{the} | \text{<S>}) P(\text{cat} | \text{the}) P(\text{ran} | \text{cat}) \\ P(\text{STOP} | \text{ran})$$

3-gram: $P(\text{the} | \text{ss}) P(\text{cat} | \text{ss the})$
 $P(\text{ran} | \text{the cat})$

Pol1

I saw the dog —
wagging, bark, in, on, jump

I saw the dog. —

IT, any word that starts a
sentence

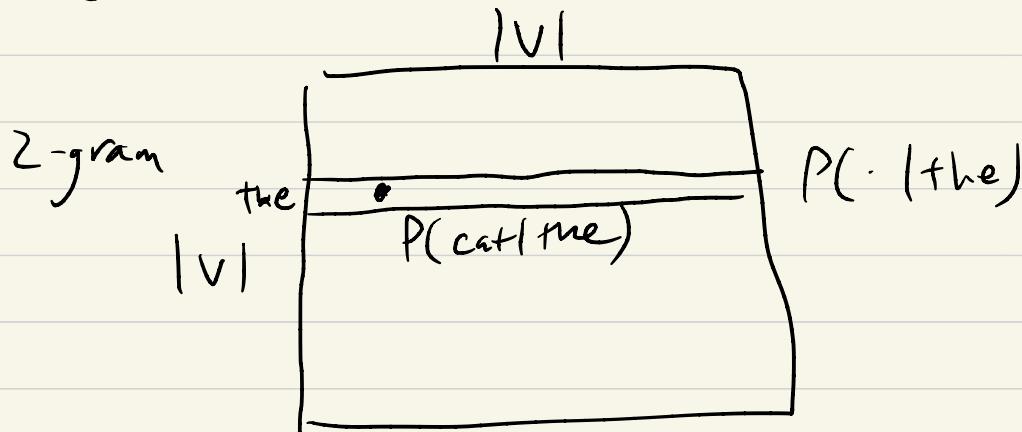
less restrictive context



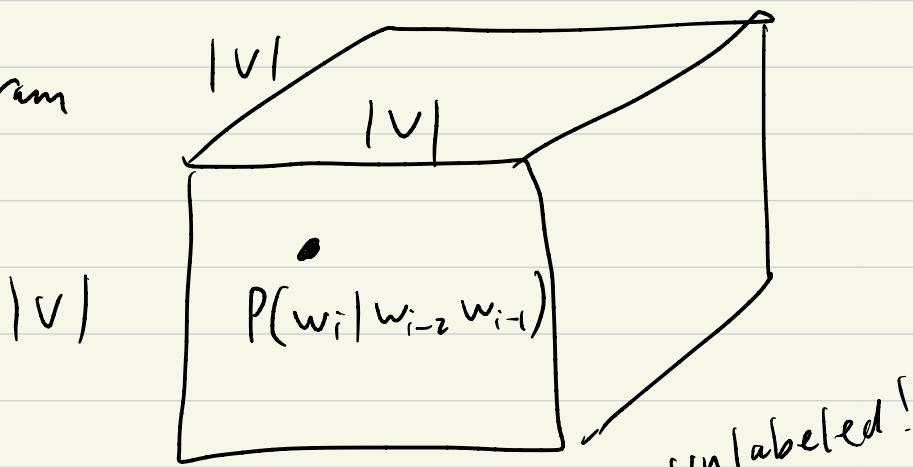
N-gram parameterization

$V = \text{vocabulary}$

Big lookup table (like HMM transitions)



3-gram



Count + normalize over a big corpus

To do well, we need $n \geq 5$
(5-gram)

I hate to go to Maui
5 words

Count ("hate to go to Maui") on the web? May be 0!

$$\Rightarrow P(\text{Maui} | \text{hate to go to}) = 0$$

\Rightarrow incorrect, should be > 0

Smoothing (in n-gram LMs)

$$P_5(w_i | w_{i-4} w_{i-3} w_{i-2} w_{i-1}) \approx \lambda P_5^{\text{raw}}(w_i | w_{i-4} \dots w_{i-1}) + (1-\lambda) P_u(w_i | w_{i-3} w_{i-2} w_{i-1})$$

raw counts

$$P_4 = \lambda P_4^{\text{raw}} + (1-\lambda) P_3$$

$$P_5 = \lambda_1 P_5^{\text{raw}} + \lambda_2 P_4^{\text{raw}} + \lambda_3 P_3^{\text{raw}} + \lambda_4 P_2^{\text{raw}} \\ + \lambda_5 P_1^{\text{raw}}$$

$$P_i^{\text{raw}} = \frac{\text{count(Marv)}}{\text{size of corpus}} > 0$$

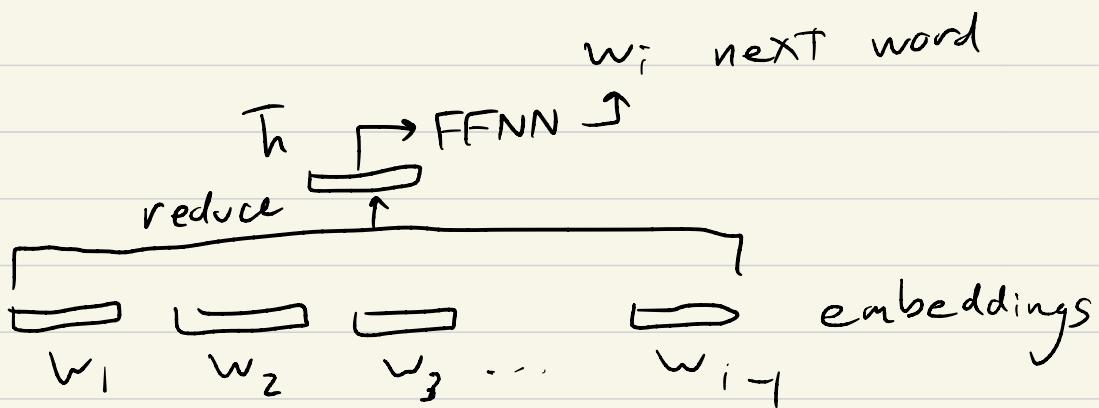
$$P_5(\text{Marv} | \dots) > 0$$

Lots of very complex tricks here!

Neural Language Models

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

$\underbrace{\quad}_{\text{all prev words, not just } n-1}$



GPT-3 reduce: Transformer

For us reduce: RNN

Also possible reduce: DAN

FFNN over $n-1$
words

d -dim vector

$$\bar{h} = \text{neural net}(w_1, \dots, w_{i-1})$$

$$P(w_i | w_1, \dots, w_{i-1}) = \text{softmax}(W \bar{h})$$

$W: |\mathcal{V}| \times d$ matrix
 weight matrix

multiclass w/ $|\mathcal{V}|$ classes

A simple ex: DAN no ordering

$$\bar{h} = \text{avg}(\bar{w}_1, \dots, \bar{w}_{i-1}) \quad d\text{-dim vector}$$

$W\bar{h}: |\mathcal{V}|$ -dim vector

Softmax: $P(w_i | \dots)$ prob dist over $|\mathcal{V}|$

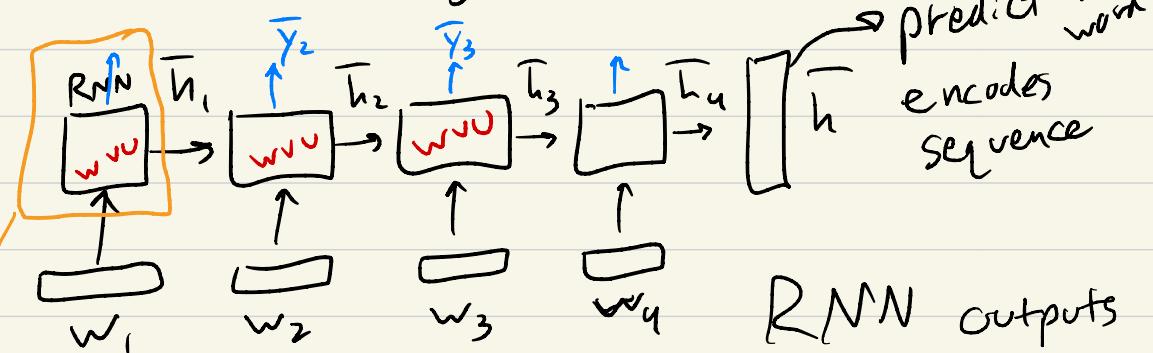
FFNN doesn't scale to large n

$$\bar{h} = \text{concat}(\bar{w}_{i-n+1}, \dots, \bar{w}_{i-1})$$

$$W: |\mathcal{V}| \times (d \cdot (n-1))$$

Concat
 $n-1$ word
 embs

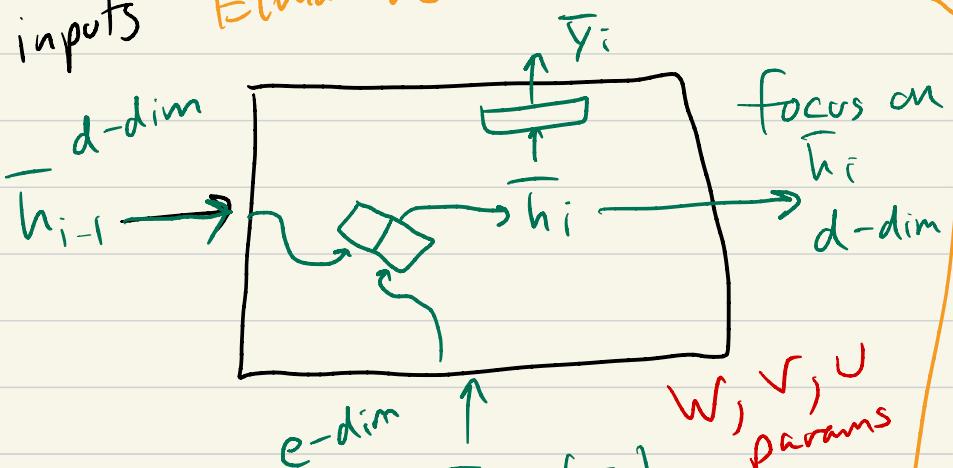
RNN encodes a sequence of vectors into a single vector "in order"



RNN outputs

$$\bar{h}_{\text{end of seq}} = \bar{h}, \bar{y}_i, \bar{h}_i \text{ at each step}$$

inputs Elman network



$$\bar{h}_i = \tanh(W_i \bar{x}_i + V \bar{h}_{i-1} + b_i)$$

$$\bar{y}_i = \tanh(U \bar{h}_i + c_i)$$

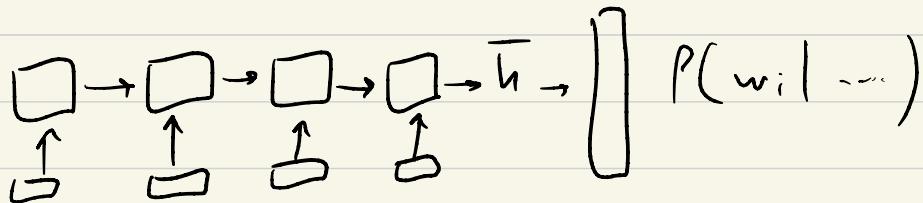
W, V, U
params

Properties

Only params are W, V, U

Params don't depend on seq. len!

Copy-paste single RNN cell



differentiable!

Compute gradients for W, V, U
w/ backprop

Form loss $(-\log P(w_i | \dots))$

Compute gradients