

CS 378 Lecture 15: Language Modeling, RNNs

Today

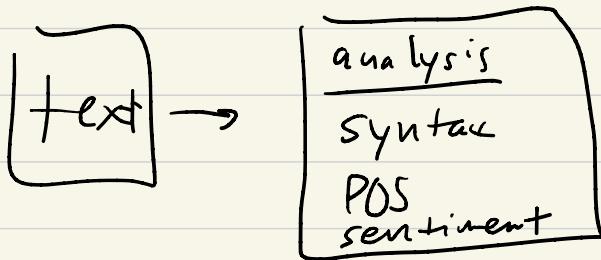
- Language modeling
- N-gram LMs
- Neural LMs
- (Start) RNNs

Announcements

- midterm
- A4 + A5
- FP custom proposals

Where we are

First half of the course : analysis
of text



input: text

output: structured analysis

same models!
(or are they?)

Next ~3 weeks:

Task: text generation ~~*~~

dialogue
summarization

input: text

output: text

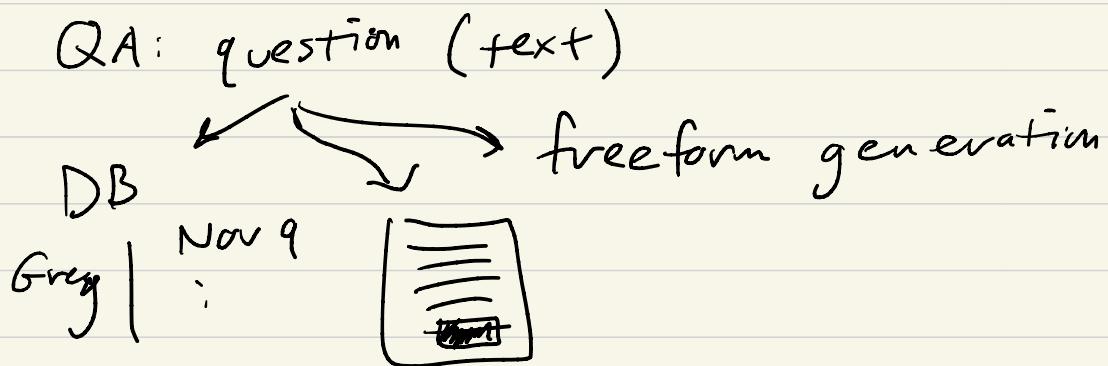
Machine translation

Language modeling: "autocomplete"
next word prediction

~~*~~ Technique: recurrent neural networks
(RNNs)
(and Transformers)

Final ~3 weeks:
Question answering] BERT,
Transformers

End of the Course: miscellaneous
applications + topics



Language Modeling

Imagine we want a distribution over
grammatical sent. in a language
 $P(\bar{w})$
PCFGs give us this

Why?

the cat ran ↙ higher prob
the the cat ran

Grammatical error correction:

I give you a sent \bar{w}
has a grammatical error

Then: $P(\bar{w})$ is low

Correct the error by finding
 \bar{w}' close to \bar{w} with
 $P(\bar{w}') >> P(\bar{w})$

Another application: machine translation

Translate \bar{w}_1 in lang 1 $\Rightarrow \bar{w}_2$ in
lang 2. Maybe we want to
make \bar{w}_2 more natural

Find \bar{w}_2' s.t. $P(\bar{w}_2') > P(\bar{w}_2)$

N-gram language models

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

True by chain rule of prob.

N-gram model: only look at the past $n-1$ words

$$P(\bar{w}) = P(w_1) P(w_2 | w_1) \cdots P(w_i | w_{i-n+1}, \dots, w_{i-1})$$

2-gram model: $\bar{w} = \text{the cat ran}$

$$P_2(\bar{w}) = P(\text{the} | \langle S \rangle) P(\text{cat} | \text{the}) - \begin{matrix} \text{start of} \\ \text{sent} \end{matrix} \cdot P(\text{van} | \text{cat}) - P(\text{STOP} | \text{van})$$

3-gram: $P(\text{the} | \langle S \rangle \langle S \rangle)$

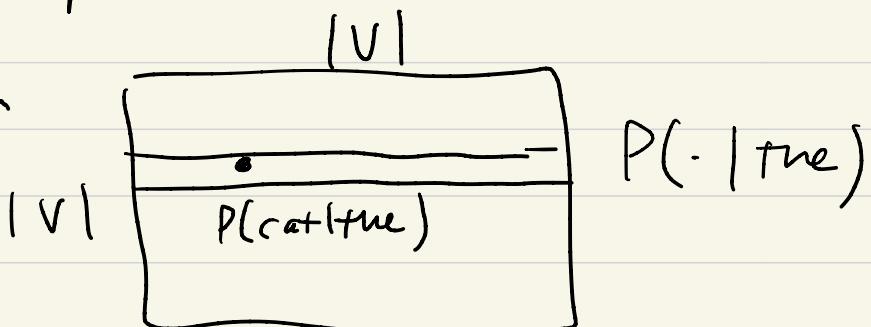
$$\cdot P(\text{cat} | \langle S \rangle \text{the})$$

$$P(\text{van} | \text{the cat}) \quad P(\text{STOP} | \dots)$$

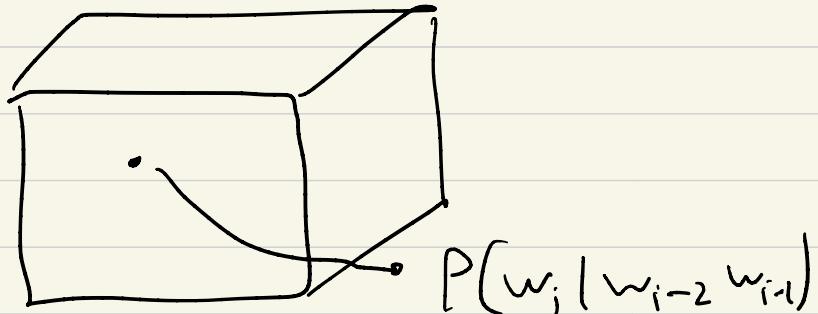
Parameters: transitions in HMM

Big lookup table

2-gram



3-gram



Get these params: count + normalize

5-gram: $P(\text{Mavil} | \text{hate to go to}) = 0$

Nobody has said this before, but
it's still reasonable!

Smoothing

Very important in n-gram LMs

$$\begin{aligned} & P(w_i | w_{i-4}, w_{i-3}, w_{i-2}, w_{i-1}) \\ \Rightarrow & \lambda_1 P^{\text{raw}}(w_i | w_{i-4}, w_{i-3}, w_{i-2}, w_{i-1}) \\ & + \lambda_2 P^{\text{raw}}(w_i | w_{i-3}, w_{i-2}, w_{i-1}) \\ & + \lambda_3 P^{\text{raw}}(w_i | w_{i-2} w_{i-1}) \quad \text{4-gram!} \\ & + \lambda_4 P^{\text{raw}}(w_i | w_{i-1}) \\ & + \lambda_5 P(w_i) \end{aligned}$$

↪ $P(\text{Mani})$

$$\lambda_1 \cdot 0 + \lambda_2 \cdot 0 + \lambda_3 \cdot 10^{-4} + \lambda_4 \cdot 10^{-2} + \dots$$

\Rightarrow reasonable prob.

Set λ carefully so this is a valid prob. dist.

Neural language models

Use a neural net to model

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

whole context, not just $n-1$

words

For now: simplify, consider

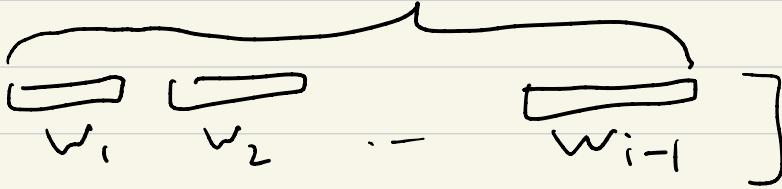
$$P(w_i | w_{i-1})$$

↑
vectors

Skip-gram:
 $P(\text{context} | \text{word})$

One idea! take context \bar{w}

$$\bar{w} \rightarrow \text{NN} \rightarrow P(w_i | w_1, \dots)$$



embed with
pre-trained
skip-gram
vectors

$P(w_i | w_1, \dots, w_{i-1})$ dist over
|V| words

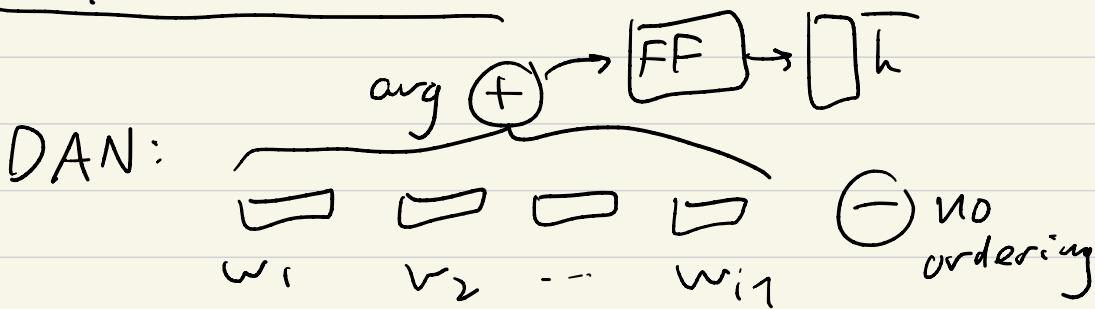
\vec{w} d-dim

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

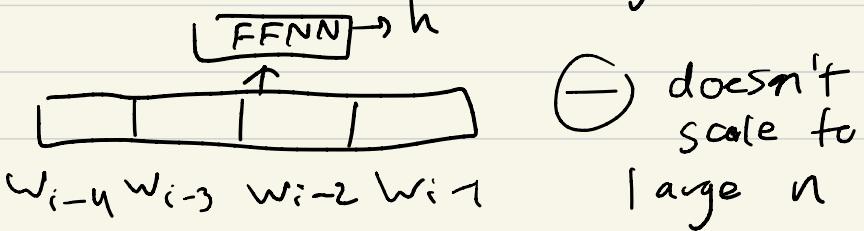
$P(w_i | w_1, \dots, w_{i-1}) = \text{softmax}(\vec{w}\vec{h})$

\vec{w} : param matrix $|V| \times d$

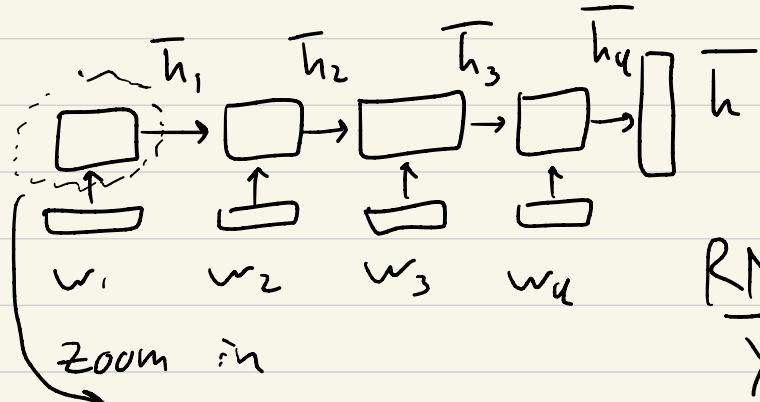
Ways to do this



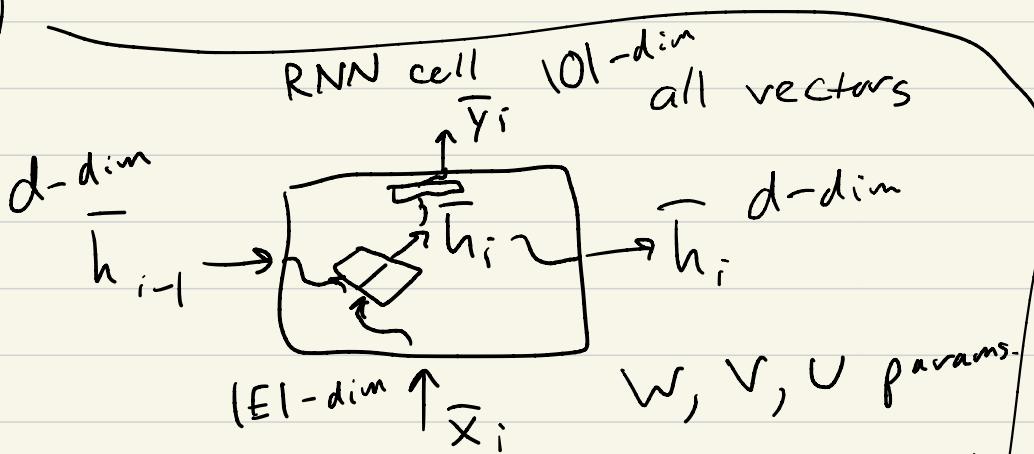
FFNN: n-1 words \vec{h} 5 gram



RNN: encodes a sequence of arbitrary length into a single vector

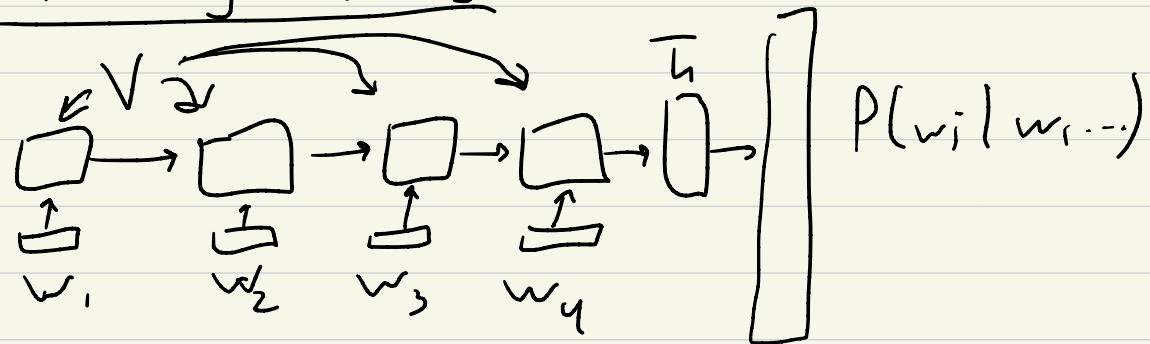


RNN outputs:
 \bar{y}_i, \bar{h}_i at each step



Elman network $\bar{h}_i = \tanh(W\bar{x}_i + V\bar{h}_{i-1})$
 $\bar{y}_i = \tanh(U\bar{h}_i)$

Training RNNs



This is a differentiable computation!

Some layers share params