

CS 378 Lecture 15: Language Modeling, n-grams, RNNs

Announcements

- Midterm back soon (Mon?)
- AY released Tues
- Vicente Ordóñez talk Fri 11am 6.302
- Custom FP proposals due in one week
(email to me)

Today

- Intro to language modeling
- N-gram LMs
- Neural LMs
- RNNs (brief!)
- Transformers on Tuesday

Recap (course so far)

Input

text $\begin{cases} \nearrow \text{classify (sentiment)} \\ \rightarrow \text{POS} \\ \searrow \text{syntactic parsing} \end{cases}$

text \rightarrow label / structure

Next few weeks:

text \rightarrow text (seq2seq)

machine translation, dialogue systems,
summarization, QA, ...

Today Language Modeling

"autocomplete" / predictive text

RNNs / Transformers

LM: a distribution $P(\bar{w})$ over sequences of words

Assign high prob. to real, grammatical, natural sentences

Grammatical error correction:

\bar{w} , try changing some function words, find \bar{w}'

$P(\bar{w}') > P(\bar{w}) \Rightarrow$ we fixed some errors!

N-gram language modeling

$\bar{w} = (w_1, \dots, w_m)$ m words

By the chain rule of probability:

$$P(\bar{w}) = P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2) \\ \text{no assumptions} \dots P(w_n | w_1 \dots w_{n-1})$$

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

2-gram LM

n-gram LM: depend on past n-1 words

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

$$n=4: w_{i-3} w_{i-2} w_{i-1}$$

the cat ran
2-gram LM: start of sequence

$$P_2(\bar{w}) = P(\text{the} | \langle s \rangle) P(\text{cat} | \text{the}) \\ P(\text{ran} | \text{cat}) P(\text{STOP} | \text{ran})$$

3-gram

$$P_3(\bar{w}) = P(\text{the} | \langle s \rangle \langle s \rangle) - \\ P(\text{cat} | \langle s \rangle \text{the}) \\ P(\text{ran} | \text{the cat})$$

Ex

I saw the dog _____ lots of completions

2⁻³

yesterday

bark

owner

sleep

store

in
at

$P(\text{Texas is})$
 $P(\text{cap. of TX})$

The capital of Texas is _____

fewer completions
trickier!
≥ 4 words

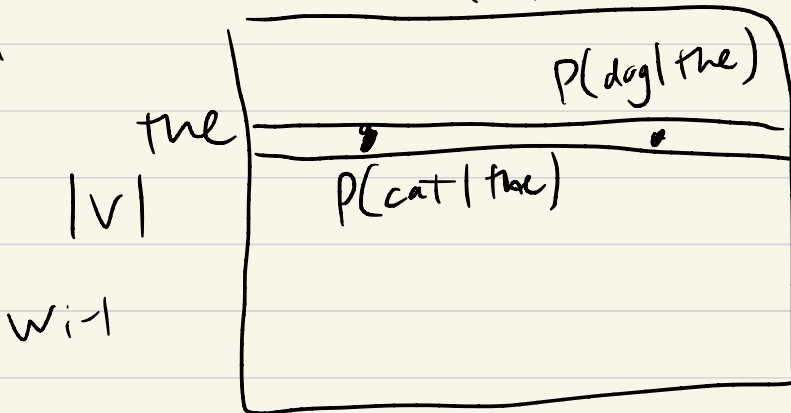
Austin
far
hot

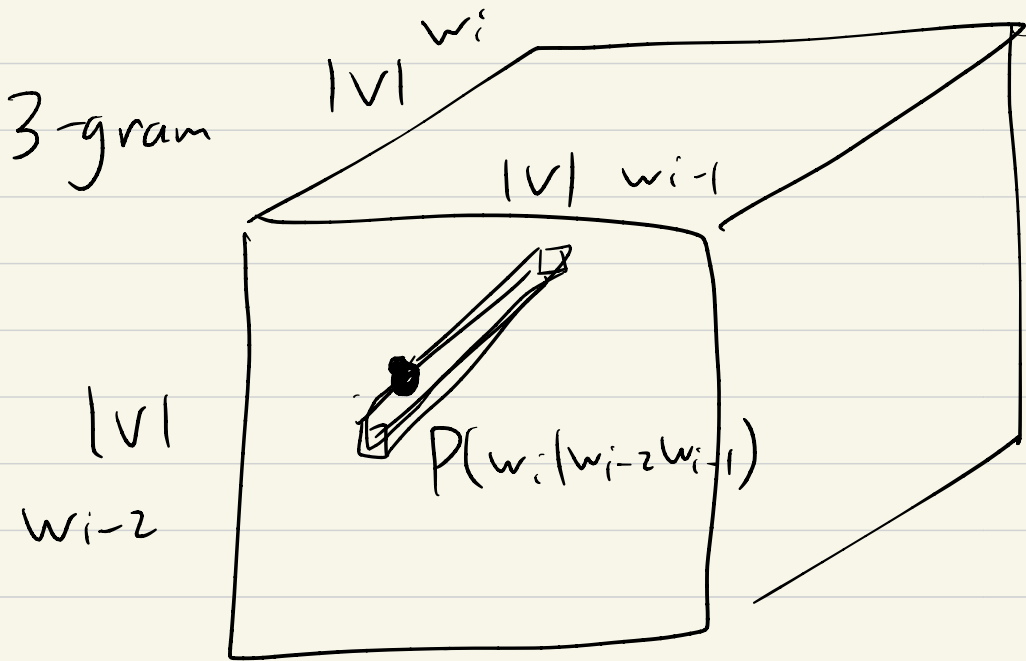
LM is a very fundamental task

≡
N-gram parameterization \checkmark vocab

Like HMM transitions: lookup table
 $|v|$ w_i

2-gram





sparse: not all 3-grams exist
 actual size $\ll |V|^3$

$n \sim 7$ is about as high as you
 can go

Parameters: Count + normalize
over a big corpus

the cat

the cat $P(\text{cat} | \text{the}) = \frac{1}{2}$

the dog

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

Ex

$n=5$, I hate to go to Maui

$\underbrace{\hspace{15em}}_{5 \text{ words}}$

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

$\frac{\text{Count}(\text{hate to go to Maui})}{\text{count}(\text{hate to go to})}$

$P_{5}^{\text{raw}} =$

$\text{count}(\text{hate to go to})$

Smoothing

$$P_5(w_i | w_{i-4} w_{i-3} w_{i-2} w_{i-1}) \\ = \lambda \cdot P_5^{\text{raw}}(w_i | w_{i-4} \text{ } -3 \text{ } -2 \text{ } -1) \\ + (1-\lambda) P_4(w_i | w_{i-3} \text{ } -2 \text{ } -1)$$

$$\lambda = 0.9$$

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

$$P_1 = P_1^{\text{raw}}(w_i) \rightarrow 0 \quad \text{any } w_i \in V \\ \text{is seen } \geq 1 \text{ time}$$

What do we store:

P_5^{raw} " |V|⁵" $P_3^{\text{raw}}, P_2^{\text{raw}}, P_1^{\text{raw}}$

P_4^{raw} " |V|⁴" tiny

Counts are sparse

distribution is not sparse

Many smoothing schemes

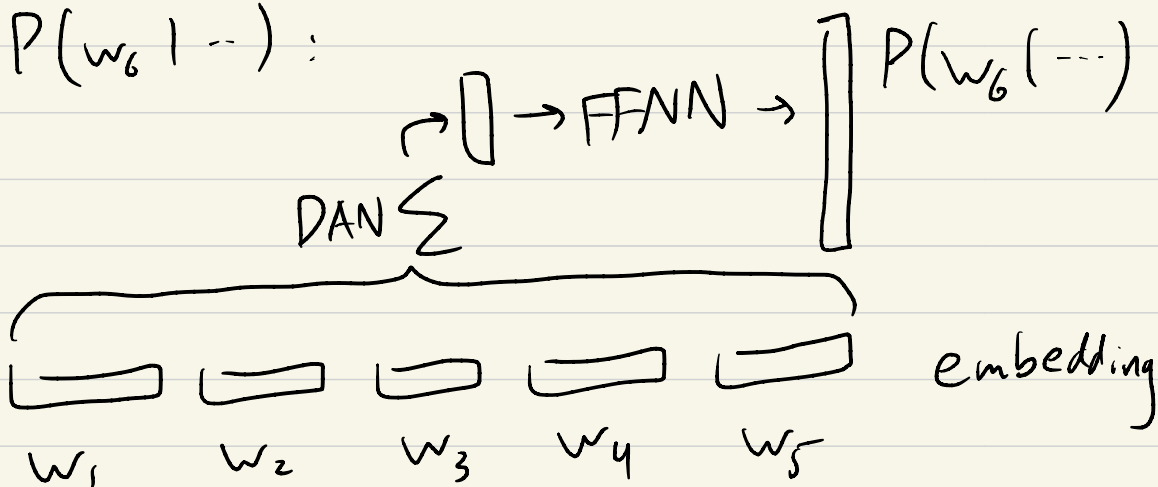
absolute discounting: sets the
 λ dynamically

Kneser - Ney smoothing

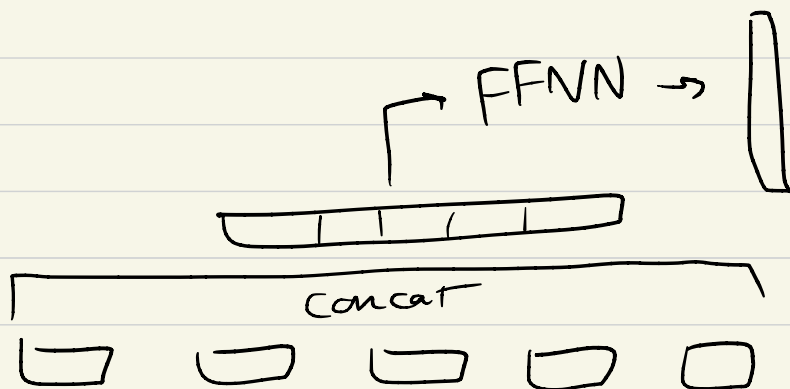
Neural Language Modeling

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

$P(w_6 | \dots)$:



Con: ignore ordering



Cons: require weights for every word in every position



the movie was —



my favorite movie —

May be hard to scale / different sizes of input

Two solutions:

- ① RNN: process sequences "uniformly"
- ② Transformer: halfway between DAN + concat

RNN will maintain a hidden state \bar{h} after processing n words

$$\bar{h}_n = \text{neural net}(w_1, \dots, w_n)$$

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

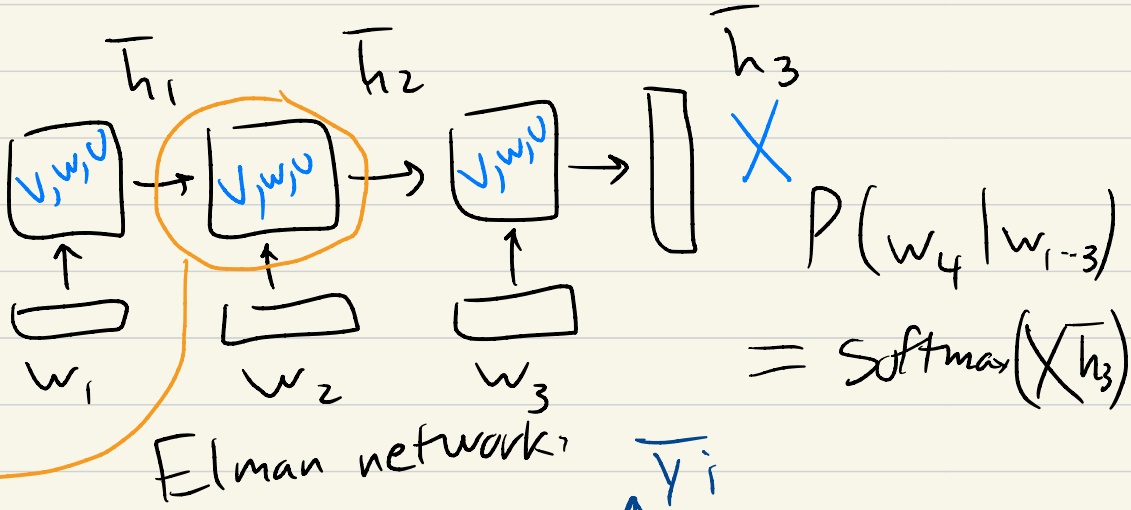
W : $|V| \times d$ matrix

\bar{h} : d -dimensional vector

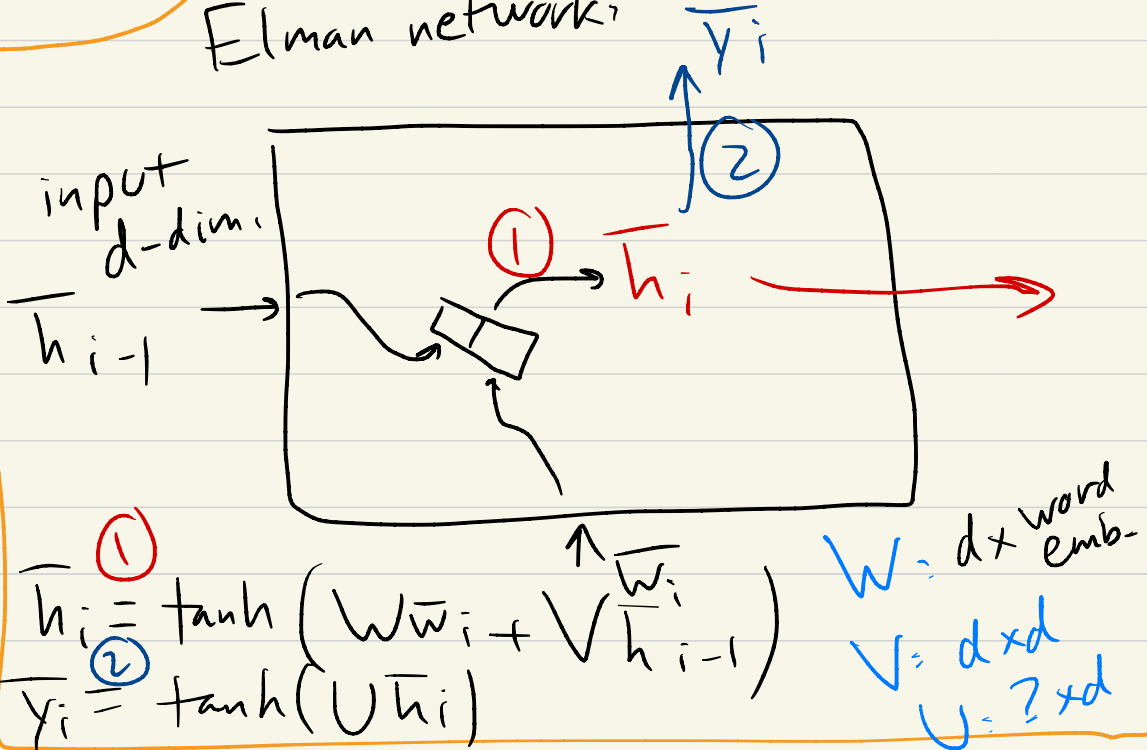
$$\bar{h}_{n+1} = \text{neural net}(\bar{h}_n, w_{n+1})$$

easy to update

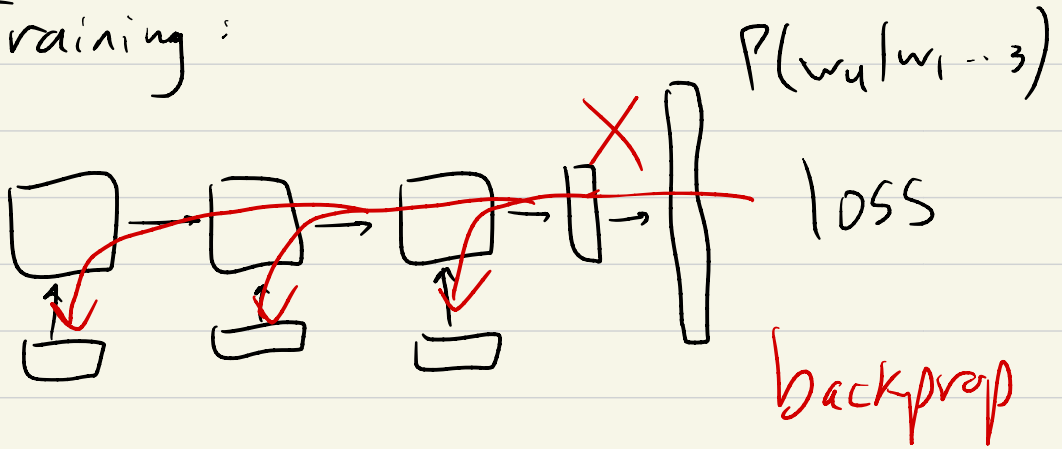
RNN encodes a sequence of vecs.
into a single vector



Elman network:



Training :



updates V, W, U

update term from every cell

Elman network

LSTM