# CS371N Lecture 10
## LM2: Self-attention, Transformers

- A2 due
- A3 out, due in 2 weeks
- Bias in embs response due

Recap Language models:
$$P(\overline{w}) = \prod_{i=1} P(w_i | w_1 \ldots w_{i-1})$$

n-gram LMs: $$P(\overline{w}) = \prod_{i=1}^{\ell} P(w_i | w_{i-n+1} \ldots w_{i-1})$$

Store probabilities explicitly (model as categorical distributions)

Estimation: count + normalize
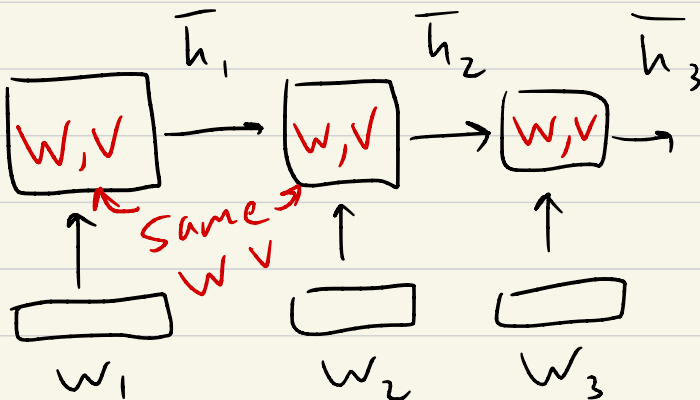
Neural LMs:

predict $w_i \mid w_1 \ldots w_{i-1}$

DANs? FFNNs?

Today

- RNNs and their shortcomings
- (Self)-attention

RNNs encode a sequence of vectors
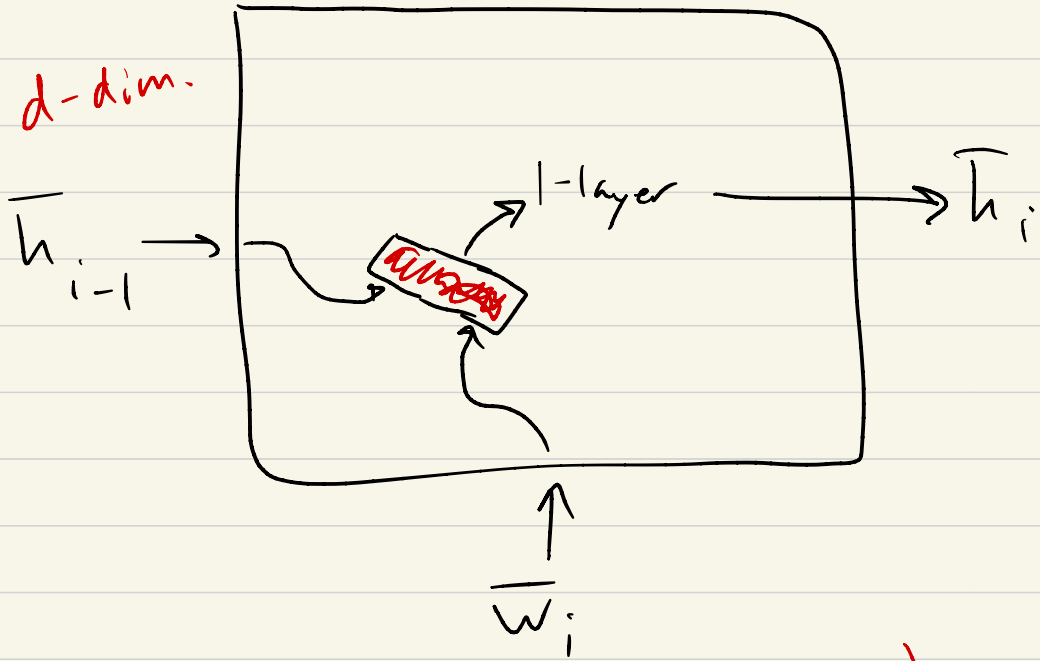  sequentially by tracking a "hidden
  state"

$\bar{h}_1 \qquad \bar{h}_2 \qquad \bar{h}_3$     "summary" of
                                                      the seq.



$P(w_4 \mid w_1 \ldots w_3)$

$= \text{softmax}$
$(U\bar{h}_3)$

RNN cell

d-dim.

$\overline{h}_{i-1} \rightarrow$

|-layer $\longrightarrow \overline{h}_i$

softmax

$\overline{w}_i$

$$\overline{h}_i = \tanh\left(W\overline{w}_i + V\overline{h}_{i-1}\right)$$
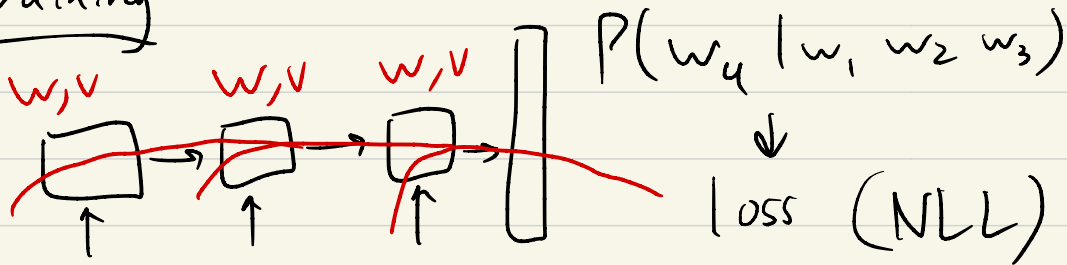
$W$: $d \times |word\ embl|$ matrix ⎤
$V$: $d \times d$ ⎦ parameters

Key property: params don't depend
  on sequence length
  $\Rightarrow$ can scale to arbitrarily
  long inputs

However, it's position-sensitive!

<u>Training</u>



$P(w_4 | w_1, w_2, w_3)$
$\downarrow$
loss (NLL)

backprop
updates to W,V accumulate over
  whole sequence

Hard to learn! $\Rightarrow$ Long short-term
                         memory net
                           (LSTM)

# Shortcomings of RNNs

"Forgetfulness": hard to track
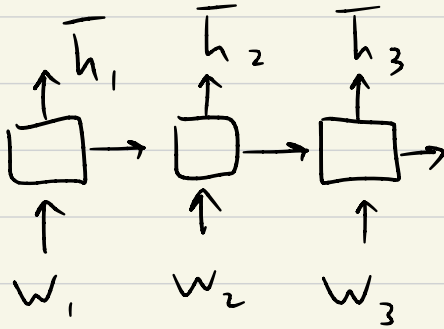information over many steps

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

$$\bar{h}_i = V\bar{h}_{i-1} \qquad h_i \approx \sum_{j=1}^{i} V^{(i-j)} W w_j$$

LSTM: "gates" to control what
parts of the vector change

$O(n)$ sequential dependence
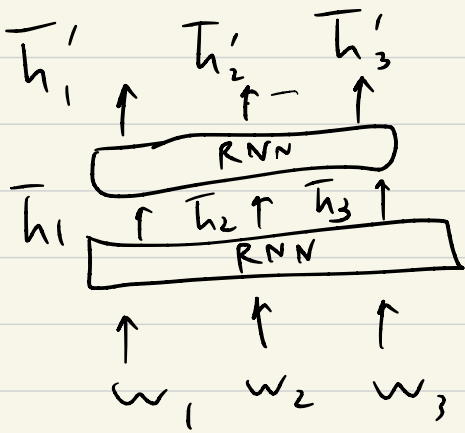
# RNN "API"



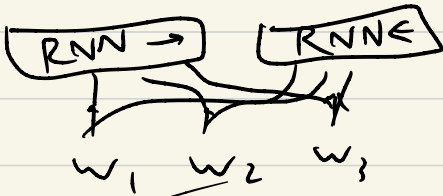$\overline{h}_3$ "context-aware" word embedding
for word $W_3$

$\overline{h}_3$ blends $\overline{h}_2$ (context) w/ $W_3$
(word 3)

RNN ( seq of vectors) $\Rightarrow$ seq of vectors
aware of context

Stack these layers

$\overline{h}'_1$    $\overline{h}'_2$    $\overline{h}'_3$

$\uparrow$    $\uparrow -$    $\uparrow$

RNN

$\overline{h}_1$   $\uparrow$   $\overline{h}_2$ $\uparrow$ $\overline{h}_3$ $\uparrow$

RNN

$\uparrow$    $\uparrow$    $\uparrow$

$w_1$   $w_2$   $w_3$

Can increase
depth

RNN $\rightarrow$    RNN $\leftarrow$

$w_1$   $w_2$   $w_3$

Transformer: layer that contextualize
words based on other words in
the sequence

$$(e'_1, e'_2, e'_3) = \text{Transformer}(e_1, e_2, e_3)$$

Running example:
Suppose we have seqs of As
and Bs of length 4
if all As → next is A
if any B → next is B

A A A A A    predict next char
A B A A B    using this sequence
B A A B B    that came before

B A A A A ... A  B   hard for RNNs
_____      to predict
     1000

<u>Attention</u> = allows us to do
"random access" on the context
to retrieve info we need

"Souped up" DAN, will add order
information next time

Keys: $K_i$ embedding $e_i$ of the sequence

query: vector representing what
we want to find

Assume key $A = \begin{bmatrix} 1 \\ 0 \end{bmatrix}^{e_A}$ $B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}^{e_B}$
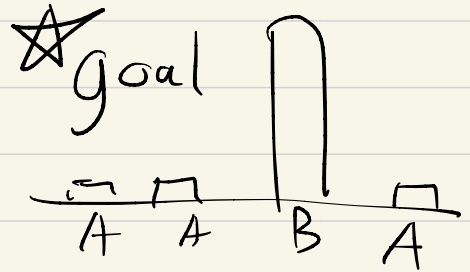(word embeddings)

$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

A    A    B    A

query: what we want to find

find Bs!          ⭐ goal

$q = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$   "B"


$\underline{\quad A \quad A \quad B \quad A \quad}$

Attention will compute a distribution over the tokens so far with higher weight for things that match $q$

Steps  ① Compute score for each key based on query

$$S_i = K_i^T q$$

$S!$  0   0   1   0     $q = \begin{bmatrix} 0 & 1 \end{bmatrix}$

$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

A    A    B    A

② Softmax scores to get prds.

$\overline{\alpha} = softmax(\overline{5})$

$\begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}$     Assume $e = 3$

$\downarrow \qquad \underbrace{\qquad\qquad} \longrightarrow \dfrac{e^1}{e^0 + e^0 + e^1 + e^0} \approx \dfrac{1}{2}$

$\begin{bmatrix} \frac{1}{6} & \frac{1}{6} & \frac{1}{2} & \frac{1}{6} \end{bmatrix}$



③ Compute the output

output $= \sum \alpha_i e_i$     weighted sum of $e_i$

$= \dfrac{1}{6}\begin{bmatrix} 1 \\ 0 \end{bmatrix} + \dfrac{1}{6}\begin{bmatrix} 1 \\ 0 \end{bmatrix} + \dfrac{1}{2}\begin{bmatrix} 0 \\ 1 \end{bmatrix} + \dfrac{1}{6}\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$= \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$

Compare to DAN:

$$\frac{1}{4}\begin{bmatrix}1\\0\end{bmatrix} + \frac{1}{4}\begin{bmatrix}1\\0\end{bmatrix} + \frac{1}{4}\begin{bmatrix}0\\1\end{bmatrix} + \frac{1}{4}\begin{bmatrix}1\\0\end{bmatrix} = \begin{Bmatrix}3/4\\1/4\end{Bmatrix}$$

$\begin{bmatrix}1/2\\1/2\end{bmatrix}$ weights the B more highly

Ideally want:

$\begin{bmatrix}1\\0\end{bmatrix}$ if all As ✓

$\begin{bmatrix}0\\1\end{bmatrix}$ if any B $\qquad \begin{bmatrix}1/2\\1/2\end{bmatrix}$ ✗

Let $q = \begin{bmatrix}0 & 10\end{bmatrix}$

softmax $\begin{bmatrix}0 & 0 & 10 & 0\end{bmatrix}$

↓

0  0  1  0   probs.

Decouple keys + queries from
  embeddings

Embedding matrix $E = \begin{matrix} A \\ A \\ B \\ A \end{matrix} \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$

target $e$

$$\underbrace{(E^T W^K)}_{} \overset{K}{\phantom{x}} (W^Q e) \qquad B$$

$$W^K = I \qquad \underset{\downarrow}{W^Q} = 10 \cdot I \qquad e = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$(\text{keep } \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}) \qquad\qquad q = \begin{bmatrix} 0 \\ 10 \end{bmatrix}$$

Parameters $W^K$ and $W^Q$ will
let us learn how to query

# Self-attention

every word is a key and query
  simultaneously

do one attention computation per
  word $\Rightarrow$ contextualized embeddings
    for each word

$(d=2)$

$E$: same embs, seq len $\times d$

$K$: same, seq len $\times d$

$Q$: seq len $\times d$  (rather than $1 \times d$)

Scores $S = Q K^T$

len×len   len×d   d×len

$S_{ij} = q_i$ (ith row of Q)

$\cdot k_j$ (jth row of k)

4 $\begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$ all pairs of qs and ks

4

For now: $K = E$

$Q = E$     $S$

word 2 $\rightarrow$ $\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \end{bmatrix}$