

N-gram language modeling:

distributions over sentences

, sequence of words

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

n-gram LM

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

$\underbrace{\quad\quad\quad}_{\text{prev. } n-1 \text{ words}}$

2-gram LM: $P(w_1 | \langle s \rangle) P(w_2 | w_1) P(w_3 | w_2) \dots$

n-gram LM \Leftrightarrow n-1 order Markov model

3-gram LM: $P(w_1 | \langle s \rangle \langle s \rangle) P(w_2 | \langle s \rangle w_1) P(w_3 | w_1, w_2) \dots$

2-gram: multinomial distribution $|V| \times |N|$ params

max likelihood estimation

Param estimation: MLE from a large corpus

$$P(\text{dog} | \text{the}) = \frac{\text{count}(\text{the, dog})}{\text{count}(\text{the})}$$

↓
table of probabilities

$$P(w | \text{the}) = \begin{cases} 0.001 & \text{house} \\ 0.005 & \text{dog} \\ 0.005 & \text{cat} \end{cases}$$

very flat dist.

why

- ① Generation: Machine translation
- ② Grammatical error correction
- ③ way to build "word2vec++"

Smoothing in n-gram LMs:

5-gram models work well

$$P(\text{Austin} \mid \text{to}) > 0 \text{ seen in data}$$

$$P(\text{Austin} \mid \text{want to go to}) = 0 \text{ if corpus isn't huge}$$

Smoothing fixes this → gives some probability to unseen instances

Absolute Discounting:

reserve mass from seen 5grams to allocate to unseen 5grams

$$P(\text{Austin} \mid \text{want to go to}) = \frac{\text{count(wtgtA)}}{\text{count(wtgt)}} + \lambda P_{AD}(A \mid \text{tgt})$$

λ is set to make this normalize
Want to go to → Maui 2 → 1.8
Want to go to → class 1 → 0.8
count = 4 → campus 1 → 0.8

(# of unique words seen in this context times k)

$$P_{AD}(A \mid \text{tgt}) = \dots + \lambda' P_{AD}(A \mid \text{tgt})$$

recursive

↳ → $P(A) > 0$ — if Austin in our data

Kneser-Ney smoothing

LM Evaluation:

regular accuracy doesn't make sense as a metric in LMs

instead use Likelihood on hold out data

$$\frac{1}{n} \sum_{i=1}^n \log P(w_i | w_1, \dots, w_{i-1}) - \text{log likelihood}$$

perplexity: $\exp(\text{avg NLL}) - \text{lower = better}$

suppose probs: $\frac{1}{4}, \frac{1}{3}, \frac{1}{4}, \frac{1}{3}$

$$\text{avg NLL (base e)} = 1.242 \quad \text{Perplexity} = 3.464 \quad \begin{matrix} \text{uses geometric mean} \\ \text{of denominators} \end{matrix}$$

ranges from 10-200 standard in LMs but not really elsewhere

Neural Language Models:

$$P(\bar{w}) = P(w_1) P(w_2 | w_1) \cdots \underbrace{P(w_n | w_1, \dots, w_{n-1})}_{\text{Model with NN}}$$

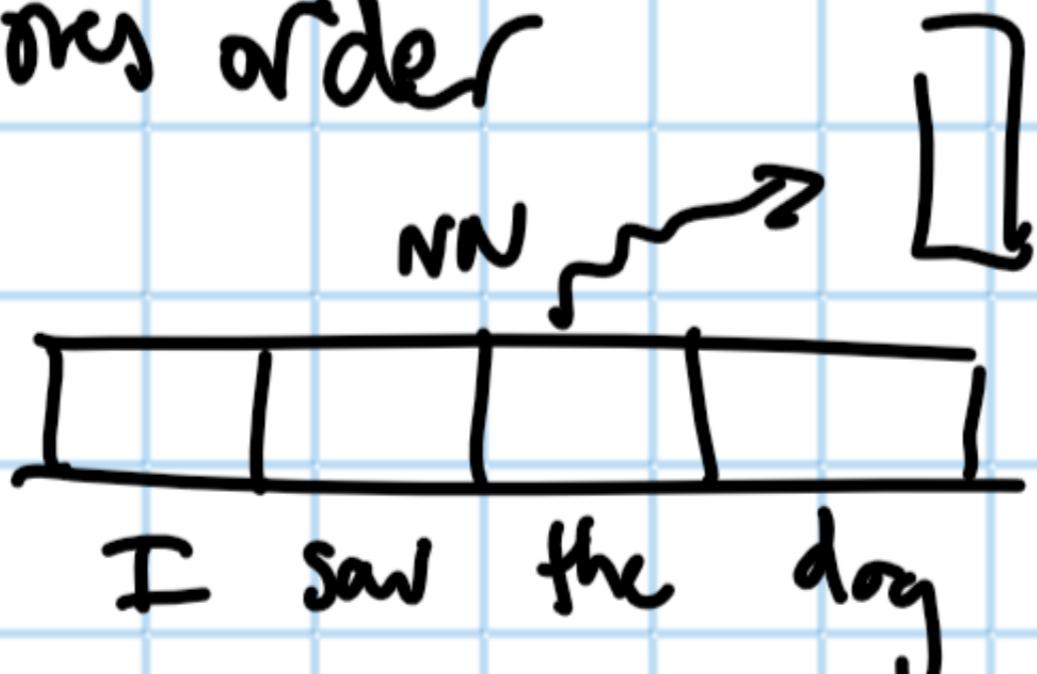
$$\text{basic: } P(w_i | w_{i-1}) = \frac{e^{v_i \cdot c_{i-1}}}{\sum_{w' \in V} e^{v_{w'} \cdot c_{i-1}}}$$

$$\text{More generally: } P(w_i | w_1, \dots, w_{i-1}) = \text{softmax}(\cup_{w_i} f(w_1, \dots, w_{i-1}))$$

$f: \text{NN to embed context}$

$f = \text{DAN} \rightarrow \text{ignores order}$

$f = \text{FFNN}$

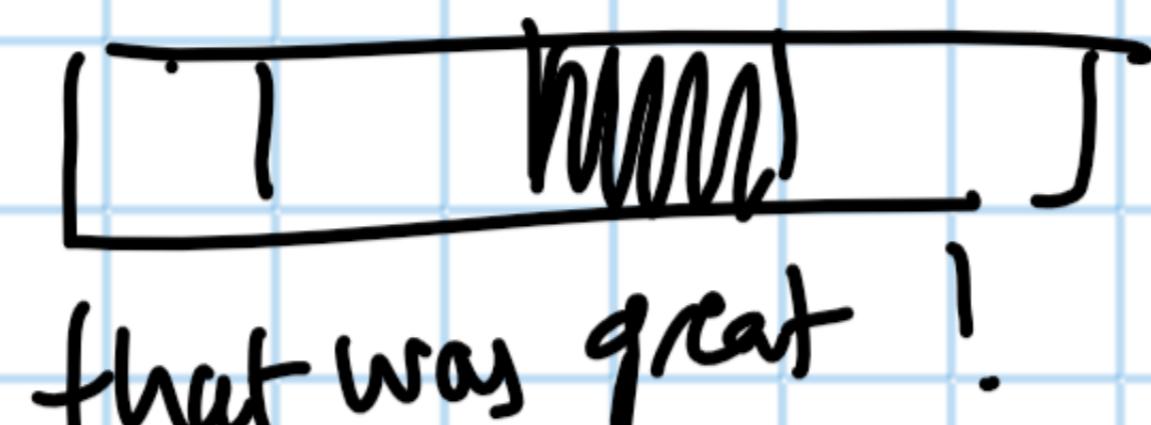
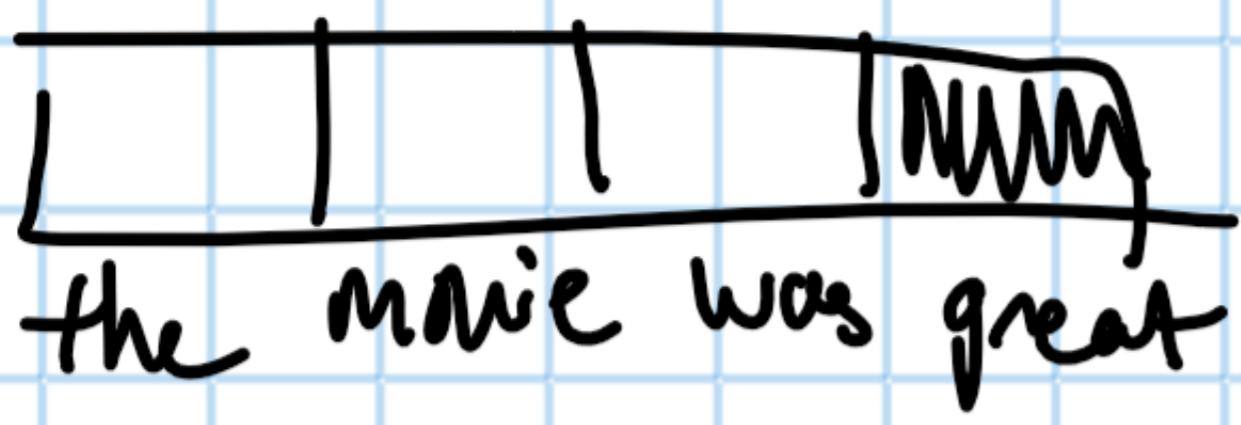


maximal softmax
 $\rightarrow P(w_i | w_1, \dots, w_{i-1})$

Mnih + Hinton (2003) tried it out but this
doesn't scale to long context (same as n-gram LM)

RNN's and their Shortcomings:

FFNN's can't handle variable length inputs : each position in the feature vector has fixed semantics



these are not related (great will be in 21 orthogonal subspaces)

different

DAN has opposite problem, where all position is gone

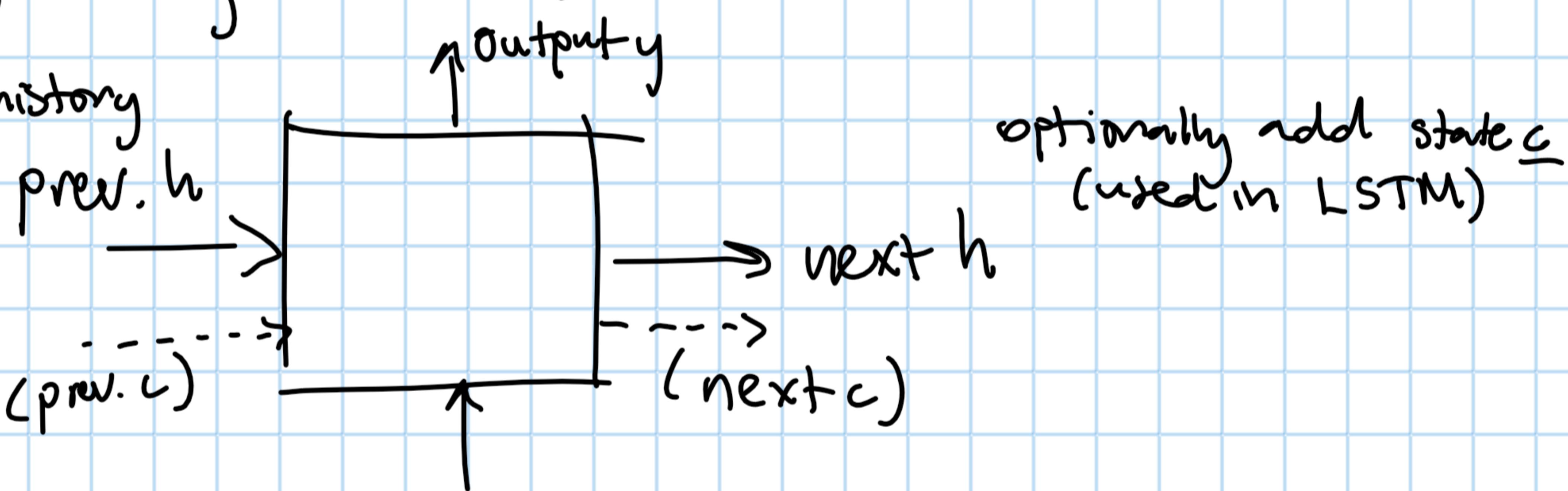
instead :

- ① process each word uniformly
- ② still using the context that the token occurs in

RNN:

cell that has input x_i , has hidden state h_i , and updates that hidden state and produces y_i (all vectors)

h represents the history



ex:
the move was great

if blue means positive sentiment,
the model recognizes great in any
part of the sentence and flips the
final output to positive

it was not great

Vanishing Gradient:

if input is at ends of spectrum, you're going to lose gradient in backpropagation

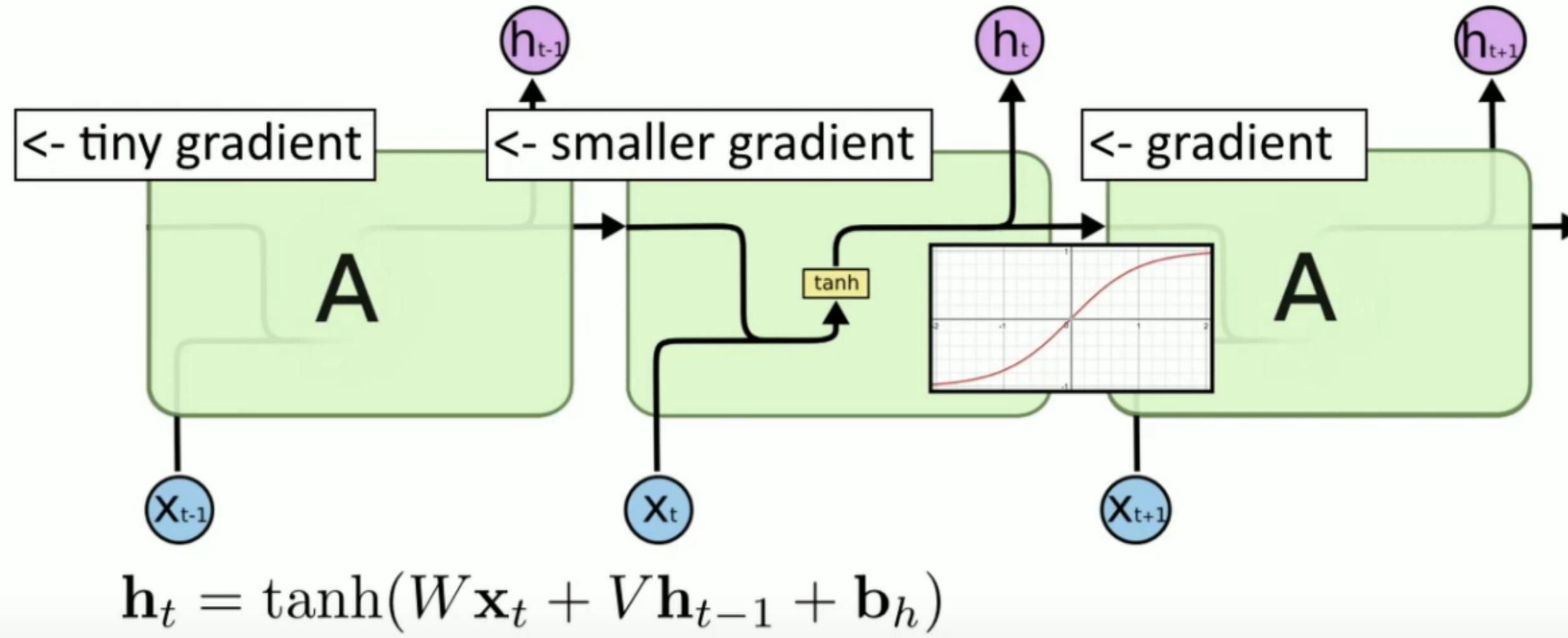
if not in $[-2, 2]$ grad is almost 0

h repeated mult by V can cause it to blow up or shrink depending on hard to get information from gradients in long scale networks

LSTM can help this, but not enough

RNNs are also slow. They don't parallelize ^{for GPU} and there are $O(n)$ non-parallel operations for encoding n items

Solution: Transformers



Attention

ex: All As: last letter is A, any B = last letter is B

AAAAA AAB
ABAAA AAB
AAAAA ABB

Attention: method to go arbitrarily far back in context from this point

, As
BAAB AAB → RNNs struggle with this (remembering context for many positions is hard)

A A B A —
keys: embeddings of seq.
Query: what we want to find
Assume: A = [1, 0] B = [0, 1]
call them e_i

Step 1: compute score for each key given query
 $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$
A A B A
0 0 1 0
Score: $s_i = k_i^T q$
Set $q = [0 \ 1]$ find B_3

Step 2: Softmax
 $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} 1/6 & 1/6 \\ 1/2 & 1/6 \end{bmatrix}$

Step 3: Compute output as weighted sum of the input

$$\text{result: } \sum_{i=1}^n \alpha_i e_i = \left(\frac{1}{6}[1 \ 0] \cdot 3 \right) + \frac{1}{2}[0 \ 1] = \left[\frac{1}{2}, \frac{1}{2} \right]$$

Compare to average: $\left[\frac{3}{4}, \frac{1}{4} \right]$

We can amplify embeddings

$$k_i = W^k e_i \quad W^k = \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix}$$

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

$$\text{softmax} \rightarrow \overline{\frac{1}{22029}}, \overline{\frac{1}{22029}}, \overline{\frac{22026}{22029}}, \overline{\frac{1}{22029}}$$

dot product attention: $s_i = k_i^T q$

scaled dot prod: $s_i = k_i^T W^Q q$

equiv to 2 wt mat: $s_i = (W^K k_i)(W^Q q)$

real attention uses 3 matrices (one for values)

Self Attention

builds on attention every word in a sequence is both a key and a query
Simultaneously

Q : seq len $\times d$ mat ($d = \text{embedding dim} = 2$ in notes)

K : seq len $\times d$ mat

$$W^Q = \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix} \quad \text{no matter value, look for Bs}$$

$$W^K = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{"booster"}$$

multiple ways to set up weights that are equiv. but this has clear structure

$$E = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$Q = EW^Q = \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$

$$K = EW^K = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}$$

SCORES: $S = QK^T$

$$S_{ij} = q_i \cdot k_j$$

$$S = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

rows represent attention scores

"I care about word 3"

Rowwise Softmax \Rightarrow distribution per row

Value mat

$$\text{output} = AS \quad (\text{actually } A(EW^V))$$

S aka attention matrix

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$$

$$V = EW^V$$

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

$n = \text{sentence length}, d = \text{hidden dim}, k = \text{kernel size}, r = \text{restricted neighborhood size}$

Quadratic complexity, but $O(1)$ sequential operations (not linear like in RNNs) and $O(1)$ "path" for words to inform each other

one head of self attention — multi heads via rand init params

produces

final mat * input = same dim as input

contextualized encoding for each word preserving length of sequence

input sequence of words \rightarrow output vector same length of sequence that represents those words in context

Multihed Self Attention

though attention can theoretically attend to multiple tokens, in practice softmax distributions become peaked



layers can help, but it would be great if we could do this in a single layer

Solution: multiple heads to do copies of attention

have 2 sets of random init matrices and compute twice

concatenate all resulting mats and multiply by a mat to bring it to output dim

Position Encodings

I visited New York and ate _____

right now, visited and ate could switch positions and we'd get the same results, unless we provide position information

encode each sequence position as an integer, add it to the word emb vector

affects input layer, so the model will learn how to attend to the 2 emb vectors

relative positional encoding: uses distance between 2 tokens

$$A_L; B_i : \text{softmax}(q_i K^T + m [-i-1, \dots, -2, -1, 0])$$

linear bias attend less for further tokens

each head will have different m

no position encodings also works with modified transformer

only looks at the past, so the model can learn position from new tokens