https://powcoder.com

Assignment Project Exam Help

Assignment Project Exam Help
Add WeChat powcoder
Sequence modeling in language
https://powcoder.com

Add WeChat powcoder

## Language modeling

► Language modeling is the task of assigning a probability to a sequence (for evaluating translation candidates).

► It is also the task of assigning probability of a word to follow a sequence of words (for language generation).

► Perfect performance in language modeling means being able to predict the next word in the sentence with a number of guesses that is fewer or less than that required by a human participant.

# Language Modeling

Formally, the language model can be formulated as:

$$P(w_{1:M}) = P(w_1)P(w_2|w_1)P(w_3|\boldsymbol{w}_{1:2})P(w_4|\boldsymbol{w}_{1:3})\cdots P(w_M|\boldsymbol{w}_{1:M-1})$$

With Markov assumption:

$$P(\boldsymbol{w}_{1:M}) \approx \prod_{m=1}^{M} P(w_m|\boldsymbol{w}_{m-n+1:m-1})$$

# Text generated by GPT-2

The fact that I work at the university makes me really nervous about my ability to find more people. That's one of the things that we're doing. My university is a very small community. Those are things that are going to have the effect, that go all these years and have the impact that it has on my life." That's one of the things that makes me really nervous about being a student. If you are a computer programmer, there are almost always somebody who you want to study. This is not how I work. Because I'm learning, I have to be very focused. If I am going to get a good job, if I want to get a career, I don't have to go out into the kitchen, you know. Now, you know, I think if I don't get the job, then I just don't have the time. Maybe I feel like I can't do much, and all of these things are going to be my own business and I want to be able to get a job. I'd certainly be looking at that. I'm at a different time. I'm going to write about my business, my business, my business and I want to talk to you about the life that I've been working for the last 30 years or so. So that's why I've been able to be here for the last 15 years through the university. My life has been extremely productive and I've had a lot of great things, so I love it."

Try it yourself: https://gpt2.ai-demo.xyz

# Noisy Channel Model for MT

- The Noisy channel model is the general framework for Statistical Machine Translation, with many variants

- Language model is the "central" component of the SMT model that is used to "select" the best translation:

$$P_{e|s}(\boldsymbol{w}^{(e)} | \boldsymbol{w}^{(s)}) \propto P_{s,e}(\boldsymbol{w}^{(e)}, \boldsymbol{w}^{(s)})$$
$$= P_{s|e}(\boldsymbol{w}^{(s)} | \boldsymbol{w}^{(e)}) \times P_e(\boldsymbol{w}^{(e)})$$

where $\boldsymbol{w}^{(e)}$ is generated from a language model, $\boldsymbol{w}^{(s)}$ is a Spanish sentence generated from a translation model $P_{s|e}(\boldsymbol{w}^{(s)} | \boldsymbol{w}^{(e)})$

# Perplexity: a metric for evaluating language models

▶ Given a text corpus of $M$ words (where $M$ could be in the millions) $w_1, w_2, w_3, \cdots w_M$, a language model function LM assigns a probability to a word based on its history.

$$\ell(\boldsymbol{w}) = \sum_{m=1}^{M} \log_2 P(w_m | w_{m-1}, w_{m-2}, \cdots, w_1)$$

▶ The perplexity of the *LM* with respect to the corpus is:

$$Perplex(\boldsymbol{w}) = 2^{-\frac{\ell(\boldsymbol{w})}{M}}$$

# What counts as a good language model?

- Good language models will assign high probability to the events in the corpus, resulting in low perplexity.
- Perplexities are corpus-specific, and perplexities for two language models are only comparable with respect to the same evaluation corpus.

# Extreme Cases of Perplexity

- In the limit of a perfect language model, probability 1 is assigned to the held-out corpus, with
$$Perplex(\boldsymbol{w}) = 2^{-\frac{1}{M} \log_2 1} = 2^0 = 1$$

- In the opposite limit, probability zero is assigned to the held-out corpus, which responds to an infinite perplexity:
$$Perplex(\boldsymbol{w}) = 2^{-\frac{1}{M} \log_2 0} = 2^\infty = \infty$$

- Assume a uniform, unigram model in which $p(w_i) = 1/V$ for all words in the vocabulary. Then

$$\log_2(\boldsymbol{w}) = \sum_{m=1}^{M} \log_2 \frac{1}{V} = - \sum_{m=1}^{M} \log_2 V = -M \log_2 V$$

$$Perplex(\boldsymbol{w}) = 2^{\frac{1}{M} M \log_2 V} = 2^{\log_2 V} = V$$

# Traditional approaches to language modeling

- Based on a $n$-order markov property
  $$P(w_{m+1}|\boldsymbol{w}_{1:m}) \approx P(w_{m+1}|\boldsymbol{w}_{m-n:m})$$
- The estimates are usually derived from corpus counts
- The role of the language model is to provide good estimates of $\hat{P}(w_{m+1}|\boldsymbol{w}_{m-n:m})$
- The maximum likelihood (MLE) estimate of $\hat{P}(w_{m+1}|\boldsymbol{w}_{m-n:m})$ is then:

$$\hat{P}_{MLE}(w_{m+1}|\boldsymbol{w}_{m-n:m}) = \frac{\#(\boldsymbol{w}_{m-n:m+1})}{\#(\boldsymbol{w}_{m-n:m})}$$

# Addressing the zero count problem

▶ Zero count for any of the n-grams resulting in zero probability for the entire corpus, meaning infinite perplexity!

▶ Add-$\alpha$ smoothing:

$$\hat{p}_{add-\alpha}(w_{m+1}|\boldsymbol{w}_{m-n:m}) = \frac{\#(\boldsymbol{w}_{m-n:m+1}) + \alpha}{\#(\boldsymbol{w}_{m-n:m}) + \alpha|V|}$$

▶ Another technique is to back off to a lower n-gram where there is a count. The Jelinek-Mercer interpolated smoothing:

$$\hat{P}_{int}(w_{m+1}|\boldsymbol{w}_{m-n:m})$$
$$= \lambda_{m-n:m}\frac{\#(\boldsymbol{w}_{m-n:m+1})}{\#(\boldsymbol{w}_{m-n:m})} + (1 - \lambda_{m-n:m})\hat{P}_{int}(\boldsymbol{w}_{m+1}|\boldsymbol{w}_{m-(n-1):m})$$

Notice this is a recursive formulation.

# Limitations of smoothed MLE based models

- ▶ Smoothing based on backoff to lower-order events, and the sequential nature of the backoff makes it hard to scale toward large n-grams.
- ▶ MLE-based language models suffer from lack of generalization across contexts.

# Neural language models

▶ Treat word prediction as a discriminative learning task, with the goal of computing the probability $P(w|u)$, where $w \in V$ is a word, and $u$ is the context that depends on *previous* words

▶ Parametrize the probability $P(w|u)$ as a function of two $K$-dimensional dense vectors, $\beta_w \in \mathbb{R}^K$ and $v_u \in \mathbb{R}^K$:

$$P(w|u) = \frac{\exp(\beta_w \cdot v_u)}{\sum_{w' \in V}(\exp(\beta_{w'} \cdot v_u))}$$

The vector of probabilities can be computed by applying the SoftMax transformation to the vector of dot products:

$$P(\cdot|u) = \text{SoftMax}([\beta_{w_1} \cdot v_u, \beta_{w_2} \cdot v_u, \cdots, \beta_{w_V} \cdot v_u])$$

▶ The word vectors $\beta_w$ are parameters of the model and can be estimated directly, e.g., using the negative log likelihood of the training corpus as the objective

# Computing the context vector

- There are different ways to compute the context vector $v$, and one effective way is to use a **Recurrent Neural Network** or RNN. The basic idea is to recurrently update the context vector while moving through a sequence.
- Let $h_m$ represent the contextual information at position $m$ in the sequence. An RNN model can be defined as:

$$x_m \triangleq \phi_{w_m}$$
$$h_m = \text{RNN}(x_m, h_{m-1})$$
$$P(w_{m+1}|w_1, w_2, \cdots, w_m) = \frac{\exp(\beta_{m+1} \cdot h_m)}{\sum_{w' \in V} \exp(\beta_{w'} \cdot h_m)}$$

where $\phi$ is a matrix of word embeddings, and $x_m$ is the word embedding for $w_m$

Sequence-to-sequence models

# Sequence-to-sequence models

- Sequence goes in, sequence comes out
- Sequence-to-sequence models are a powerful learning framework that have found success in a wide range of applications
  - Automatic Speech Recognition (ASR): sound stream goes in, text comes out
  - Machine Translation (MT): source language sentence goes in, target language sentence comes out
  - Image captioning: Image goes in, caption comes out
  - text summarization: whole text goes in, summary comes out
  - Automatic email responses: Generating automatic responses to incoming emails
  - etc. etc.

# The encoder decoder architecture

▶ The encoder network converts the source sentence into a vector or a matrix representation; the decoder network then converts the encoding into a sentence in the target language

$$z = \text{Encode}(w^{(s)})$$
$$w^{(t)} | w^{(s)} \sim \text{Decode}(z)$$

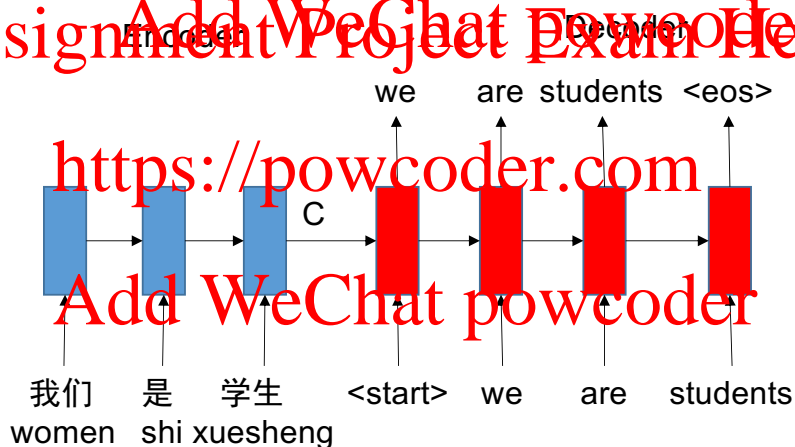where the second line means the decoder defines the conditional probability $P(w^{(t)}|w^{(s)})$.

▶ The decoder is typically a recurrent neural network (e.g., LSTM) that generates one word at a time, while recurrently updating a hidden state.

▶ The encoder decoder networks are trained end-to-end from parallel sentences.

Encoder decoder

Encoder          Decoder

we    are  students  <eos>

C

我们    是    学生      <start>    we      are      students
women  shi xuesheng

# Training objective

If the output layer of the decoder is a logistic function, then the entire network can be trained to maximize the conditional log-likelihood (or minimize the negative log-likelihood):

$$\log P(\mathbf{w}^{(t)}|\mathbf{w}^{(s)}) = \sum_{m=1}^{M^{(t)}} P(w_m^{(t)}|\mathbf{w}_{1:m-1}^{(t)}, \mathbf{z})$$

$$w_m^{(t)}|\mathbf{w}_{1:m-1}^{(t)}, \mathbf{w}^{(s)} \sim \text{SoftMax}\left(\boldsymbol{\beta} \cdot \mathbf{h}_{m-1}^{(t)}\right)$$

where $\mathbf{h}_{m-1}^{(t)}$ is a recurrent function of the previously generated text $\mathbf{w}_{1:m-1}^{(t)}$ and the ecoding $\mathbf{z}$, and $\boldsymbol{\beta} \in \mathbb{R}^{(V^{(t)} \times K)}$ is the matrix of output word vectors for the $V^{(t)}$ words in the target language vocabulary

# The LSTM variant

▶ In the LSTM variant of the encoder network, the encoder is set to the final hidden state of an LSTM on the source sentence:

$$\boldsymbol{h}_m^{(s)} = \text{LSTM}(\boldsymbol{x}_m^{(s)}, \boldsymbol{h}_{m-1}^{(s)})$$

$$\boldsymbol{z} \triangleq \boldsymbol{h}_{M^{(s)}}^{(s)}$$

where $\boldsymbol{x}^{(s)}$ is the embedding of the source language word $w_m^{(s)}$.

▶ The encoding then provides the initial hidden state for the decoder LSTM:

$$\boldsymbol{h}_0^{(t)} = \boldsymbol{z}$$

$$\boldsymbol{h}_m^{(t)} = \text{LSTM}(\boldsymbol{x}_m^{(t)}, \boldsymbol{h}_{m-1}^{(t)})$$

where $\boldsymbol{x}_m^{(t)}$ is the embedding of the target language word $w_m^{(t)}$
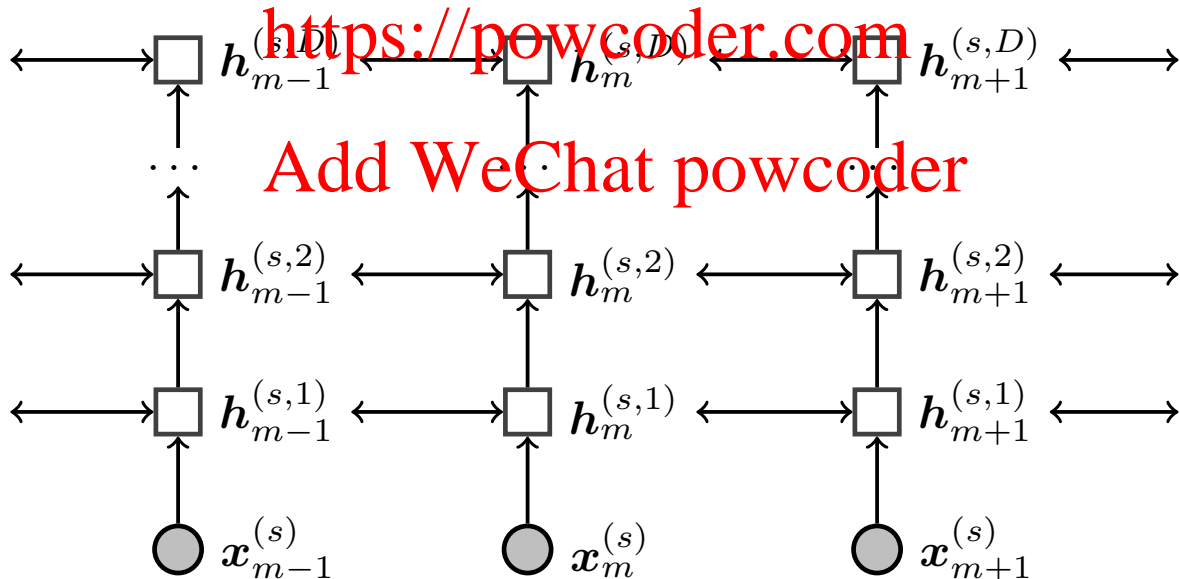
# Tweaking the encoder decoder network

- ▶ Adding **layers**: The encoder and decoder network can be implemented as **deep LSTMs** with multiple layers of hidden state
- ▶ Adding **attention** to give more weight to particular word or words in the source language when generating a word in the target language

# Multi-layered LSTMs

Each hidden state $\boldsymbol{h}^{(s,i)}$ at layer $i$ is treated as the input to an LSTM at layer $i+1$:

$$\boldsymbol{h}_m^{(s,1)} = \text{LSTM}(\boldsymbol{x}_m^{(s)}, \boldsymbol{h}_{m-1}^{(s)})$$

$$\boldsymbol{h}_m^{(s,i+1)} = \text{LSTM}(\boldsymbol{h}_m^{(s,i)}, \boldsymbol{h}_{m-1}^{(s,i+1)}), \forall i \geq 1$$

# Neural attention

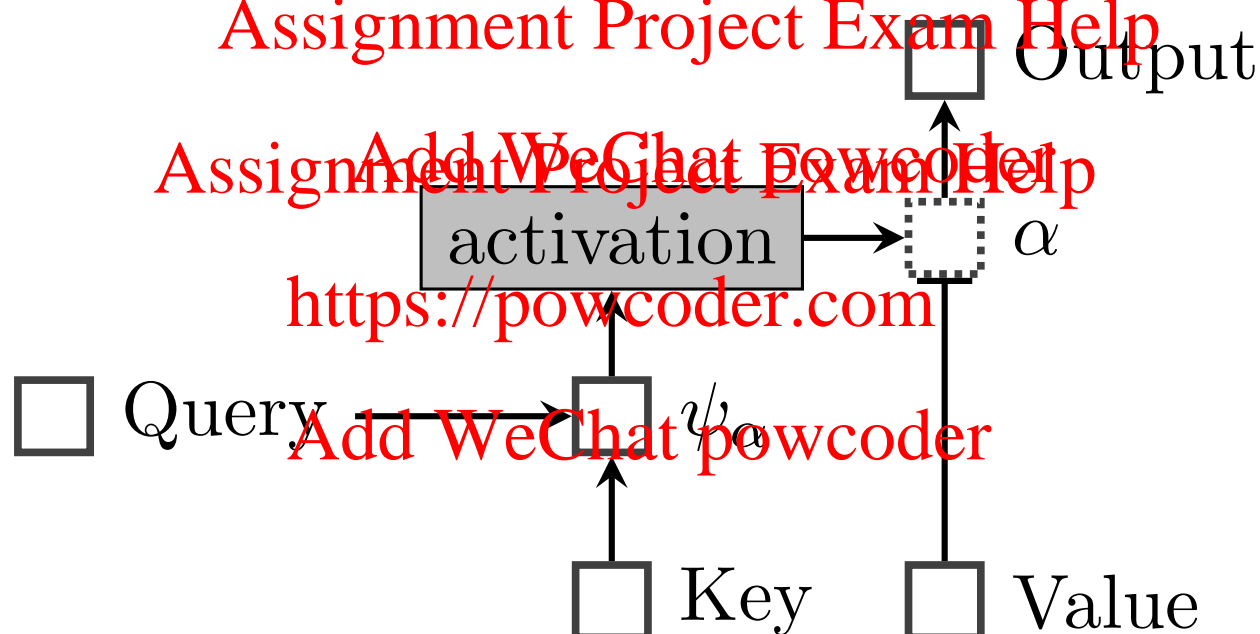- Attention can be thought of as using a query to select from a memory of key-value pairs, with the keys, values and queries all being vectors

- For each key $n$ in the memory, we compute a score with respect to the query $m$, which measures the "comptability" between the key and the query

- The scores are passed through an activation function, typically softmax, which results in a vector of non-negative numbers of length $N$, which equal to the size of the memory:
  $$[\alpha_{m \to 1}, \alpha_{m \to 2}, \cdots, \alpha_{m \to N}]$$

- Multiply each value in the memory $v_n$ by the attention $\alpha_{m \to n}$, and sum them up, we get the output of the attention.

- The attention is typical concatenated with the decoding hidden state to output the target word

"Querying"

activation $\longrightarrow$ $\alpha$

Output

Query $\longrightarrow$ $\psi_\alpha$

Key

Value

## Step by step computation of attention

▶ Computing compatibility with a two-layer feedforward network:

$$\psi_\alpha(m, n) = v_\alpha \cdot \tanh(\Theta_\alpha[h_m^{(t)}; h_n^{(s)}])$$

▶ Softmax attention

$$\alpha_{m \to n} = \frac{\exp \psi_\alpha(m, n)}{\sum_{n'=1}^{M^{(s)}} \exp \psi_\alpha(m, n')}$$

▶ Compute the weighted average over columns of $Z$

$$c_m = \sum_{n=1}^{M^{(s)}} \alpha_{m \to n} z_n$$

▶ incorporate the context vector into the decoding model:

$$\tilde{h}_m^{(t)} = \tanh\left(\Theta_c[h_m^{(t)}, c_m]\right)$$

$$P(w_{m+1}^{(t)} | w_{1:m}^{(t)}, w^{(s)} \propto \exp\left(\beta_{w_{m+1}^{(t)}} \cdot \tilde{h}_m^{(t)}\right)$$

# Seq2seq: Initialization

▶ Word embeddings

$$\boldsymbol{E} = embed\left(\begin{bmatrix} women \\ shi \\ xuesheng \\ START \\ we \\ are \\ students \\ EOS \end{bmatrix}\right) = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & -0.5 \\ -0.5 & 0.5 \\ 0 & 0 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \\ 0.1 & 0.2 \\ 1 & 1 \end{bmatrix}$$

▶ RNN parameters

$$\boldsymbol{W}^{(s)} = \begin{bmatrix} 0.3 & 0 \\ 0 & 0.3 \end{bmatrix} \quad \boldsymbol{U}^{(s)} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} \quad \boldsymbol{b}^{(s)} = \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$$

$$\boldsymbol{W}^{(t)} = \begin{bmatrix} 0.3 & 0 \\ 0 & 0.3 \end{bmatrix} \quad \boldsymbol{U}^{(t)} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} \quad \boldsymbol{b}^{(t)} = \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}$$

## Encoder

$$z \triangleq h^{(s)}$$

$$h_n^{(s)} = \tanh(W^{(s)} \times x + U^{(s)} \times h + b)$$

$$h_1^{(s)} = \tanh(W^{(s)} \times E[women] + U^{(s)} \times 0 + b^{(s)}) = \begin{bmatrix} 0.3364 \\ 0.5717 \end{bmatrix}$$

$$h_2^{(s)} = \tanh(W^{(s)} \times E[shi] + U^{(s)} \times h_1^{(s)} + b^{(s)}) = \begin{bmatrix} 0.3153 \\ 0.7289 \end{bmatrix}$$

$$h_3^{(s)} = \tanh(W^{(s)} \times E[xuesheng] + U^{(s)} \times h_2^{(s)} + b^{(s)}) = \begin{bmatrix} 0.0086 \\ 0.5432 \end{bmatrix}$$

$$C = h_3^{(s)} = \begin{bmatrix} 0.0086 \\ 0.5432 \end{bmatrix}$$

where $C$ is a context vector

Decoder

$$h_m^{(t)} = \tanh(W^{(t)} \times x + U^{(t)} \times h + b)$$

$$h_1^{(t)} = \tanh(W^{(t)} \times E[\text{START}] + U^{(t)} \times h + b^{(t)}) = \begin{bmatrix} 0.6929 \\ 0.2482 \end{bmatrix}$$

$$h_2^{(t)} = \tanh(W^{(t)} \times E[\text{we}] + U^{(t)} \times h_1^{(t)} + b^{(t)}) = \begin{bmatrix} 0.6203 \\ 0.2123 \end{bmatrix}$$

$$h_3^{(t)} = \tanh(W^{(t)} \times E[\text{are}] + U^{(t)} \times h_2^{(t)} + b^{(t)}) = \begin{bmatrix} 0.6039 \\ 0.1870 \end{bmatrix}$$

$$h_4^{(t)} = \tanh(W^{(t)} \times E[\text{students}] + U^{(t)} \times h_3^{(t)} + b^{(t)}) = \begin{bmatrix} 0.6220 \\ 0.0980 \end{bmatrix}$$

# Softmax over similarities between hidden layers and target embeddings

$$score_1\left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix}\right) = softmax\left(\begin{bmatrix} \boldsymbol{h}_1^{(t)} \times \boldsymbol{E}[we] \\ \boldsymbol{h}_1^{(t)} \times \boldsymbol{E}[are] \\ \boldsymbol{h}_1^{(t)} \times \boldsymbol{E}[students] \\ \boldsymbol{h}_1^{(t)} \times \boldsymbol{E}[EOS] \end{bmatrix}\right) = \begin{bmatrix} 0.1913 \\ 0.2000 \\ 0.1733 \\ 0.4354 \end{bmatrix}$$

$$score_2\left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix}\right) = softmax\left(\begin{bmatrix} \boldsymbol{h}_2^{(t)} \times \boldsymbol{E}[we] \\ \boldsymbol{h}_2^{(t)} \times \boldsymbol{E}[are] \\ \boldsymbol{h}_2^{(t)} \times \boldsymbol{E}[students] \\ \boldsymbol{h}_2^{(t)} \times \boldsymbol{E}[EOS] \end{bmatrix}\right) = \begin{bmatrix} 0.1999 \\ 0.2070 \\ 0.1826 \\ 0.4116 \end{bmatrix}$$

$$score_3\left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix}\right) = softmax\left(\begin{bmatrix} \boldsymbol{h}_3^{(t)} \times \boldsymbol{E}[we] \\ \boldsymbol{h}_3^{(t)} \times \boldsymbol{E}[are] \\ \boldsymbol{h}_3^{(t)} \times \boldsymbol{E}[students] \\ \boldsymbol{h}_3^{(t)} \times \boldsymbol{E}[EOS] \end{bmatrix}\right) = \begin{bmatrix} 0.2012 \\ 0.2098 \\ 0.1867 \\ 0.4023 \end{bmatrix}$$

$$score_4\left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix}\right) = softmax\left(\begin{bmatrix} \boldsymbol{h}_4^{(t)} \times \boldsymbol{E}[we] \\ \boldsymbol{h}_4^{(t)} \times \boldsymbol{E}[are] \\ \boldsymbol{h}_4^{(t)} \times \boldsymbol{E}[students] \\ \boldsymbol{h}_4^{(t)} \times \boldsymbol{E}[EOS] \end{bmatrix}\right) = \begin{bmatrix} 0.2037 \\ 0.2147 \\ 0.1959 \\ 0.3857 \end{bmatrix}$$

$$P(Y|X) = score_1 \times score_2 \times score_3 \times score_4$$

# Attention

▶ The idea: Different context vectors are used when generating target words at different time steps, e.g.,

$$C_1 = 0.98 \times \boldsymbol{h}_1^{(s)} + 0.01 \times \boldsymbol{h}_2^{(s)} + 0.01 \times \boldsymbol{h}_3^{(s)}$$

$$C_2 = 0.01 \times \boldsymbol{h}_1^{(s)} + 0.98 \times \boldsymbol{h}_2^{(s)} + 0.01 \times \boldsymbol{h}_3^{(s)}$$

$$C_3 = 0.01 \times \boldsymbol{h}_1^{(s)} + 0.01 \times \boldsymbol{h}_2^{(s)} + 0.98 \times \boldsymbol{h}_3^{(s)}$$

$$C_m = \sum_n \alpha_{m \to n} \times \boldsymbol{h}_n^{(s)}$$

$$\alpha_{m \to n} = \frac{\exp(score(\boldsymbol{h}_m^{(t)}, \boldsymbol{h}_n^{(s)}))}{\sum_{n'=1} \exp(score(\boldsymbol{h}_m^{(t)}, \boldsymbol{h}_{n'}^{(s)}))}$$

$$score(\boldsymbol{h}_m^{(t)}, \boldsymbol{h}_n^{(s)}) = {\boldsymbol{h}_m^{(t)}}^\top \boldsymbol{h}_n^{(s)}$$

▶ Other scoring variants exist

# Computing attention

are

$h'_2$

$c_2$

$a_{2,3}$

$a_{2,2}$

$a_{2,1}$

$h_1^{src}$  $h_2^{src}$  $h_3^{src}$  $h_2^{tgt}$

我们  是  学生  <start>  we  are  students

women  shi xuesheng

## Other attention variants

▶ Additive attention:

$$\psi_\alpha(m, n) = \boldsymbol{v}_\alpha \cdot \tanh(\boldsymbol{\Theta}_\alpha[\boldsymbol{h}_m^{(t)} + \boldsymbol{h}_n^{(s)}])$$

▶ Multiplicative attention

$$\psi_\alpha(m, n) = \boldsymbol{h}_m^{(t)^\top} \boldsymbol{\Theta}_\alpha \boldsymbol{h}_n^{(s)}$$

# Drawbacks of RNNs

- ► For RNNs, input is processed as a sequence. The computation of each state ($h_i$) depends on the previous state $h_{i-1}$.
- ► This prevents parallel computation for all tokens in the input sequence simultaneously, making it difficult to take full advantage of modern computation architecture to increase speed.
- ► We can imagine a network in which each token in the sequence interacts with any other token in the sequence. Conceptually, this can be viewed as a fully-connected graph where each token is a node in the graph, and the computation of its hidden state depends on all other tokens in the graph.
- ► With this approach, the computation of the hidden state of a hidden state $h_i$ does not depends on the computation of another hidden state. It only depends on the input sequence.
- ► With this approach, the order information would have to be captured separately, with position encoding.