

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help
Sequence-to-sequence models

<https://powcoder.com>

Add WeChat powcoder

Sequence-to-sequence models

<https://powcoder.com>

- ▶ 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.

Translation



Translation

<https://powcoder.com>



The encoder decoder architecture

<https://powcoder.com>

- ▶ 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

Add WeChat Exam Help
 $z = \text{ENCODE}(w^{(s)})$

$w^{(t)} | w^{(s)} \sim \text{DECODE}(z)$

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

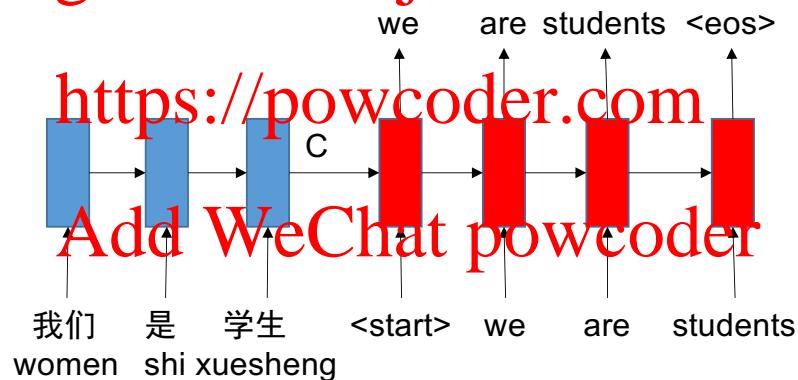
- ▶ The decoder can be a recurrent neural network (e.g., LSTM) that generates one word at a time, while recurrently updating a hidden state, or it can be a transformer.
- ▶ The encoder decoder networks are trained end-to-end from parallel sentences.

Encoder decoder

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help



Training objective

<https://powcoder.com>

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):

[Assignment Project Exam Help](https://powcoder.com)

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

$$w_m^{(t)} | \mathbf{w}_{1:m-1}^{(t)}, \mathbf{w}^{(s)}, \text{SoftMax}(\beta \cdot h_{m-1}^{(t)})$$

where $\mathbf{h}_{m-1}^{(t)}$ is a recurrent function of the previously generated text $\mathbf{w}_{1:m-1}^{(t)}$ and the encoding \mathbf{z} , and $\beta \in \mathbb{R}^{(V^{(t)} \times K)}$ 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:
Assignment Project Exam Help

Assignment Project Exam Help

$$z \triangleq h_{M^{(s)}}^{(s)}$$

https://powcoder.com

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

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

$$h_0^{(t)} = z$$

$$h_m^{(t)} = \text{LSTM}(x_m^{(t)}, h_{m-1}^{(t)})$$

where $x_m^{(t)}$ is the embedding of the target language word $w_m^{(t)}$

Tweaking the encoder decoder network

<https://powcoder.com>

Assignment Project Exam Help

- ▶ Adding layers: The encoder and decoder network can be implemented as **deep LSTMs** with multiple layers of hidden state <https://powcoder.com>
- ▶ Adding **attention** to give more weight to particular word or words in the source language when generating a word in the target language [Add WeChat powcoder](https://powcoder.com)

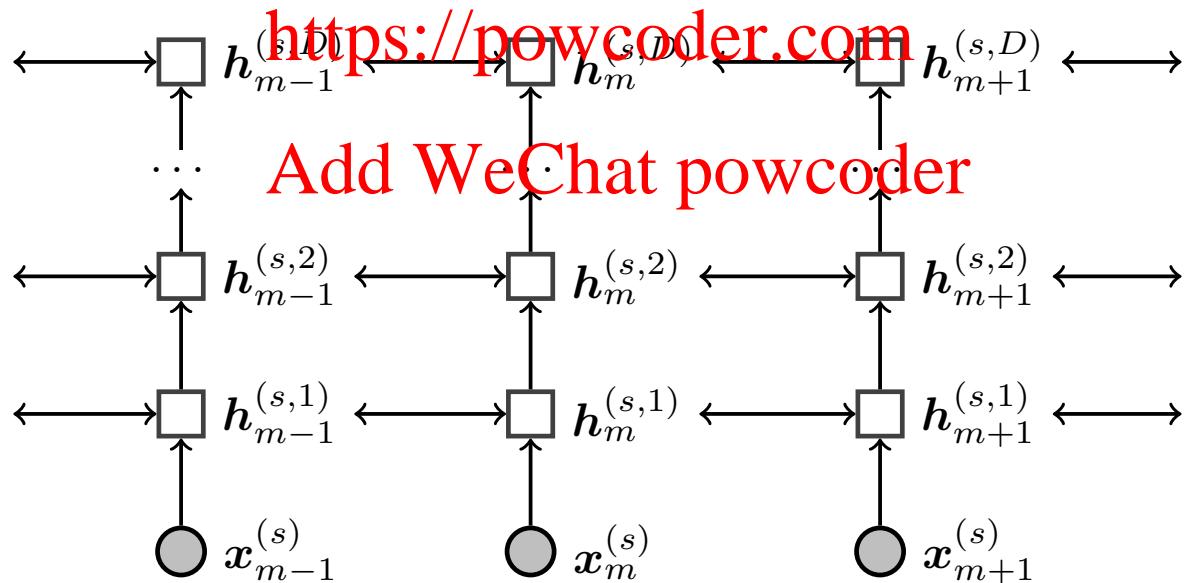
Multi-layered LSTMs

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

$$h_m^{(s,1)} = \text{LSTM}(x_m^{(s)}, h_{m-1}^{(s)})$$

$$h_m^{(s,i+1)} = \text{LSTM}(h_m^{(s,i)}, h_{m-1}^{(s,i+1)}), \forall i \geq 1$$

Add WeChat powcoder



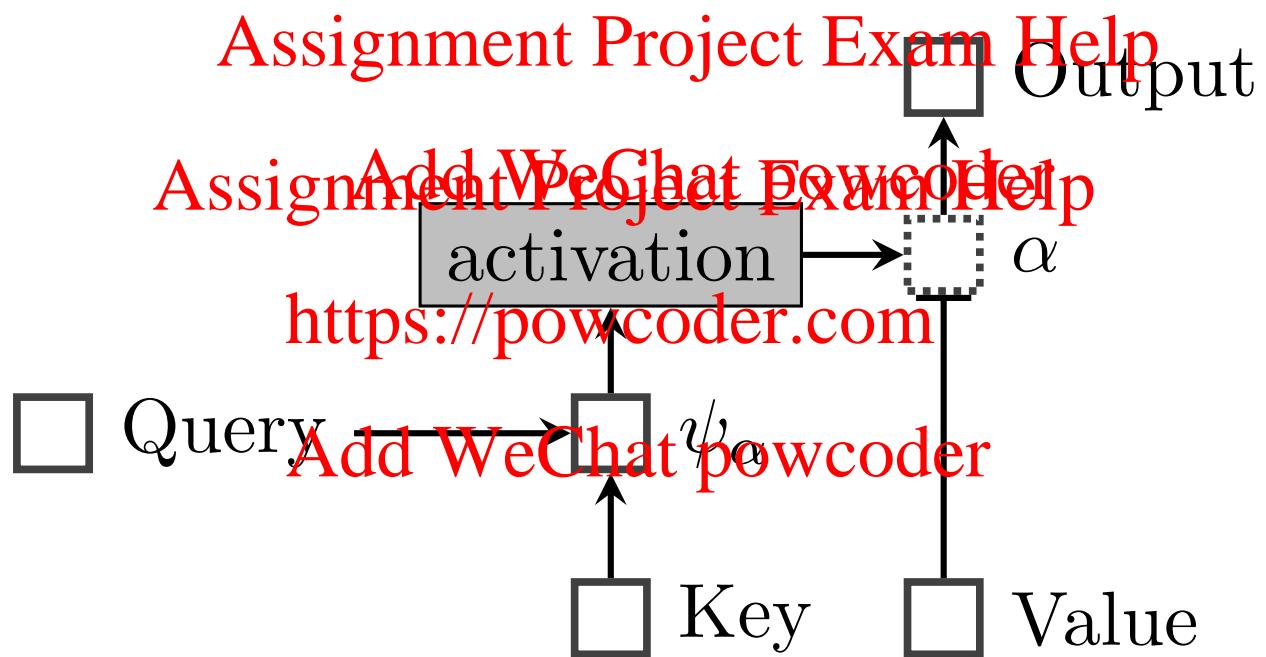
Neural attention

<https://powcoder.com>

- ▶ 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 "compatibility" 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:
$$[\alpha_{m \rightarrow 1}, \alpha_{m \rightarrow 2}, \dots, \alpha_{m \rightarrow N}]$$
- ▶ Multiply each value in the memory v_n by the attention $\alpha_{m \rightarrow n}$, and sum them up, we get the output of the attention.
- ▶ The attention is typically concatenated with the decoding hidden state to output the target word

“Querying”

<https://powcoder.com>



Step by step computation of attention

- ▶ Computing compatibility with a two-layer feedforward network:
<https://powcoder.com>

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

Assignment Project Exam Help

- ▶ Softmax attention

Add WeChat powcoder

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

<https://powcoder.com>

- ▶ Compute the weighted average over columns of \mathbf{Z}

Add WeChat powcoder

$$\mathbf{c}_m = \sum_{n=1}^{M^{(s)}} \alpha_{m \rightarrow n} \mathbf{z}_n$$

- ▶ incorporate the context vector into the decoding model:

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

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

Seq2seq: Initialization

- Word embeddings <https://powcoder.com>

$$\text{Assignment Project Exam Help}$$
$$\text{Assignment Add WeChat Exam Help}$$
$$\mathbf{E} = \text{embed} \quad \begin{bmatrix} \text{women} \\ \text{shi} \\ \text{xuesheng} \\ \text{START} \\ \text{we} \\ \text{are} \\ \text{students} \\ \text{EOS} \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & -0.5 \\ -0.5 & 0.5 \\ 0 & 0 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \\ 0.1 & 0.2 \\ -1 & 1 \end{bmatrix}$$
$$\text{Add WeChat powcoder}$$

<https://powcoder.com>

- RNN parameters

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

Encoder

<https://powcoder.com>

$$z \triangleq h^{(s)} \quad \text{Assignment Project Exam Help}$$

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

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

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

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

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

where C is a context vector

Decoder

<https://powcoder.com>

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

$$\mathbf{h}_1^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[we] + \mathbf{U}^{(t)} \times \mathbf{h}_m^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6929 \\ 0.2482 \end{bmatrix}$$

$$\mathbf{h}_2^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[we] + \mathbf{U}^{(t)} \times \mathbf{h}_1^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6203 \\ 0.2123 \end{bmatrix}$$

$$\mathbf{h}_3^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[are] + \mathbf{U}^{(t)} \times \mathbf{h}_2^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6039 \\ 0.1870 \end{bmatrix}$$

$$\mathbf{h}_4^{(t)} = \tanh(\mathbf{W}^{(t)} \times \mathbf{E}[students] + \mathbf{U}^{(t)} \times \mathbf{h}_3^{(t)} + \mathbf{b}^{(t)}) = \begin{bmatrix} 0.6220 \\ 0.0980 \end{bmatrix}$$

Softmax over similarities between hidden layers and target embeddings <https://powcoder.com>

$$score_1 \left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_1^{(t)} \times E[we] \\ h_1^{(t)} \times E[are] \\ h_1^{(t)} \times E[students] \\ h_1^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.1913 \\ 0.2000 \\ 0.1733 \\ 0.4354 \end{bmatrix}$$

$$score_2 \left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_2^{(t)} \times E[we] \\ h_2^{(t)} \times E[are] \\ h_2^{(t)} \times E[students] \\ h_2^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.1999 \\ 0.2070 \\ 0.1826 \\ 0.4116 \end{bmatrix}$$

$$score_3 \left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_3^{(t)} \times E[we] \\ h_3^{(t)} \times E[are] \\ h_3^{(t)} \times E[students] \\ h_3^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.2012 \\ 0.2098 \\ 0.1867 \\ 0.4023 \end{bmatrix}$$

$$score_4 \left(\begin{bmatrix} we \\ are \\ students \\ EOS \end{bmatrix} \right) = softmax \left(\begin{bmatrix} h_4^{(t)} \times E[we] \\ h_4^{(t)} \times E[are] \\ h_4^{(t)} \times E[students] \\ h_4^{(t)} \times E[EOS] \end{bmatrix} \right) = \begin{bmatrix} 0.2037 \\ 0.2147 \\ 0.1959 \\ 0.3857 \end{bmatrix}$$

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

Attention

- The idea: Different context vector C is produced when generating target words at different time steps, e.g.,

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

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

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

<https://powcoder.com>

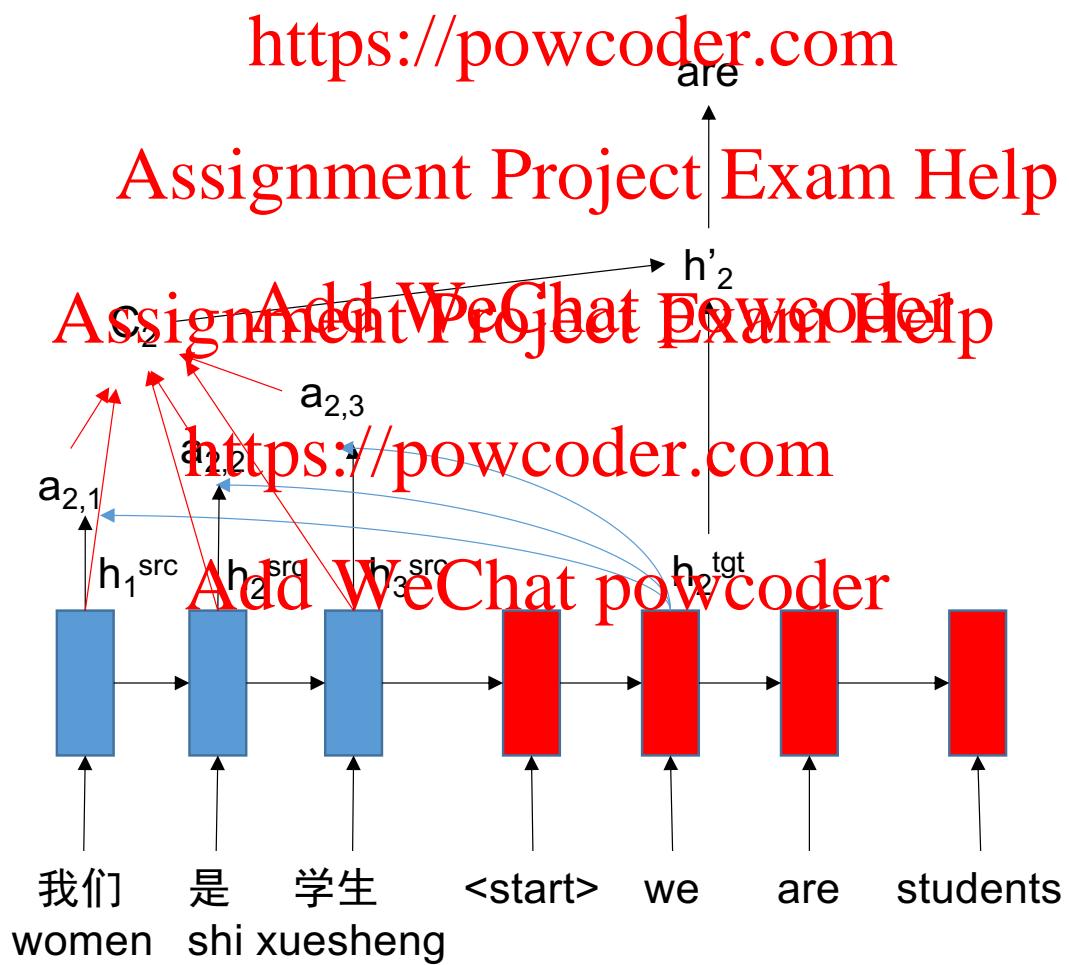
$$C_m = \sum_n \alpha_{m \rightarrow n} \times \mathbf{h}_n^{(s)}$$

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

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

- Other scoring variants exist

Computing attention



Other attention variants

<https://powcoder.com>

Assignment Project Exam Help

- ▶ Additive attention:

Add WeChat powcoder
 $\psi_\alpha(m, n) = \mathbf{v}_\alpha \cdot \tanh(\Theta_\alpha[\mathbf{h}_m^{(t)} + \mathbf{h}_n^{(s)}])$

<https://powcoder.com>

- ▶ Multiplicative attention

Add WeChat powcoder

$$\psi_\alpha(m, n) = \mathbf{h}_m^{(t)^\top} \Theta_\alpha \mathbf{h}_n^{(s)}$$

Drawbacks of RNNs

- ▶ For RNNs, input is processed as a sequence. The computation of each state (\mathbf{h}_i) depends on the previous state \mathbf{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 \mathbf{h}_i does not depend 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.