

Natural Language Processing

Yue Zhang Westlake University







Chapter 16

Working With Two Texts

Contents



- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network

Contents

VestlakeNLP

• 16.1 Sequence to sequence models

- 16.1.1 Model 1: Seq2seq using LSTM
- 16.1.2 Model 2: Adding more source features using target-tosource attention
- 16.1.3 Model 3: Copying form the source
- 16.1.4 Subword Embedding
- 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network



- Using neural models, hand-crafted features can be replaced by Dense Representation.
- A pipelined approach with sub-systems can be replaced by a simple end-to-end system.





- Using neural models, hand-crafted features can be replaced by Dense Representation.
- A pipelined approach with sub-systems can be replaced by a simple end-to-end system.





• Take an example of Machine Translation:



VestlakeNLP

- Sequence-to-sequence is useful for *more than just MT*
- Many NLP tasks can be phrased as sequence-to-sequence:
 - Machine translation(source text \rightarrow target text)
 - Summarization (long text \rightarrow short text)
 - **Dialogue** (previous utterances \rightarrow next utterance)

Encoder-Decoder Structure

VestlakeNLP

- Encoder for learning a neural input representation.
- **Decoder** for constructing the output sequence.

For Machine Translation:



Encoder-Decoder Structure

VestlakeNLP

- Encoder for learning a neural input representation.
- **Decoder** for constructing the output sequence.



For Dialogue System:

Contents

VestlakeNLP

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network

WestlakeNLP

• Encoding and Decoding with LSTM

Given an input sentence: $X_{1:n} = x_1, x_2, ..., x_n$



Generating an output sentence: $Y_{1:m} = y_1, y_2, ..., y_m$

WestlakeNLP

• Encoder: Bi-directional LSTM



$$\vec{\mathbf{h}}_{i}^{enc} = LSTM\Big(\vec{\mathbf{h}}_{i-1}^{enc}, emb(x_{i}) \\ \overleftarrow{\mathbf{h}}_{i}^{enc} = LSTM\Big(\overleftarrow{\mathbf{h}}_{i+1}^{enc}, emb(x_{i}) \\ \mathbf{h}^{enc} = [\overrightarrow{\mathbf{h}}_{n}^{enc}; \overleftarrow{\mathbf{h}}_{1}^{enc}] \\ \text{To decoder}$$



- Encoder: some tips
 - A useful practice to add pseudo start < *s* > and end tokens </*s* > to the beginning and end of the input.
 - It is useful to use $\sum_{i=1}^{n} emb(x_i)$ as the start hidden states $\vec{\mathbf{h}}_{0}^{enc}$ and $\overleftarrow{\mathbf{h}}_{n+1}^{enc}$ in both directions instead of zeros.

WestlakeNLP

• **Decoder**: output of encoder is used as an initial state h_0^{dec}



$$\mathbf{h}_{1}^{dec} = LSTM\Big(\mathbf{h}_{0}^{dec}, emb'(\langle s \rangle)\Big)$$

Predicting the first word *y*_{*i*}:

$$\mathbf{o}_1 = \mathbf{W} \mathbf{h}_1^{dec}$$

 $\mathbf{p}_1 = softmax(\mathbf{o}_1) \ \mathbf{p}_1 \in \mathbb{R}^{|V|}$

 $P(y_1|X_{1:n}) = \mathbf{p}_1[y_1]$

VestlakeNLP

• Decoder: each subsequent step *i*



$$\mathbf{h}_{i}^{dec} = LSTM\left(\mathbf{h}_{i-1}^{dec}, emb'(y_{i-1})\right)$$
$$\mathbf{o}_{i} = \mathbf{W}\mathbf{h}_{i}^{dec}$$
$$\mathbf{p}_{i} = softmax(\mathbf{o}_{i})$$

WestlakeNLP

• Decoder: each subsequent step *i*



$$\mathbf{h}_{i}^{dec} = LSTM\left(\mathbf{h}_{i-1}^{dec}, emb'(y_{i-1})\right)$$
$$\mathbf{o}_{i} = \mathbf{W}\mathbf{h}_{i}^{dec}$$
$$\mathbf{p}_{i} = softmax(\mathbf{o}_{i})$$

A greedy local Greedy search:

1

$$y_{i} = argmax_{y_{i'}} P(y_{i'}|X_{1:n}, Y_{1:i-1})$$
¹⁷

WestlakeNLP

• Training:

Given an input sentence: $X_{1:n} = x_1, x_2, ..., x_n$

$$L = -\sum_{i=1}^{N} \sum_{j=1}^{m_{i}} \log \left(P(y_{j}^{i} | X_{i}, Y_{i:j-1}^{i}) \right)$$

Generating an output sentence: $Y_{1:m} = y_1, y_2, ..., y_m$



- Target vocabulary
 - OOV: Vocabulary, in the scale of 10⁴ is very large. So highly infrequent words are not considered to make the model smaller.
 - **<UNK>**: adding a special token **<**UNK> to address OOV problem, so that gradients can be computed.
 - Disallowing the <UNK> from being generated may lead to inconsistency between training and testing. So some postprocessing can be conducted.

Contents



- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network



- Lack of model 1
 - Model 1 gives competitive results for MT and summarization, however, source representation vector h^{enc} is the only connection point.
 - RNNs face challenges in maintaining long-range dependencies between target and source words, when sentences are too long.





- Intuition beyond model 1
 - Keep the source at hand in each generating step.
 - Consider the hidden states for each source word simultaneously as the representation of whole source sentence rather than using only the last hidden states.



WestlakeNLP

Construct a context vector using attention mechanism



VestlakeNLP

• Construct a context vector using attention mechanism



$$\mathbf{c}_{i} = \sum_{j=1}^{n} \alpha_{ij} \mathbf{h}_{j}^{enc}$$
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{n} \exp(e_{ik})}$$
$$e_{ij} = \mathbf{v}_{a}^{T} \tanh(\mathbf{W}_{a} \mathbf{h}_{i-1}^{dec} + \mathbf{U}_{a} \mathbf{h}_{j}^{enc})$$

WestlakeNLP

Construct a context vector using attention mechanism



Introducing *c*^{*i*} into LSTM gates respectively:

$$\begin{aligned} \mathbf{u}_{i} &= \sigma(\mathbf{W}_{u}emb'(y_{i-1}) + \mathbf{U}_{u}\mathbf{h}_{i-1}^{dec} + \mathbf{C}_{u}\mathbf{c}_{i}) \\ \mathbf{f}_{i} &= \sigma(\mathbf{W}_{f}emb'(y_{i-1}) + \mathbf{U}_{f}\mathbf{h}_{i-1}^{dec} + \mathbf{C}_{f}\mathbf{c}_{i}) \\ \mathbf{o}_{i} &= \sigma(\mathbf{W}_{o}emb'(y_{i-1}) + \mathbf{U}_{o}\mathbf{h}_{i-1}^{dec} + \mathbf{C}_{o}\mathbf{c}_{i}) \\ \tilde{\mathbf{c}}_{i}^{dec} &= \mathbf{u}_{i} \circ \tanh(\mathbf{W}_{h}emb'(y_{i-1}) + \mathbf{U}_{h}\mathbf{h}_{i-1}^{dec} + \mathbf{C}_{h}\mathbf{c}_{i}) + \mathbf{f}_{i} \circ \tilde{\mathbf{c}}_{i-1}^{dec} \\ \mathbf{h}_{i}^{dec} &= \mathbf{o}_{i} \circ \tilde{\mathbf{c}}_{i}^{dec} \end{aligned}$$

WestlakeNLP

• Generation using attention mechanism as model 1



$$\mathbf{h}_{i}^{dec} = LSTM\left(\mathbf{h}_{i-1}^{dec}, emb'(y_{i-1})\right)$$
$$\mathbf{o}_{i}' = \mathbf{W}\mathbf{h}_{i}^{dec}$$
$$\mathbf{p}_{i} = softmax(\mathbf{o}_{i}')$$



- Beam search:
 - Denote the probability of an output sequence as model 1:

 $P(Y_{1:i}|X_{1:n}) = P(y_1|X_{1:n})P(y_2|y_1, X_{1:n}) \dots P(y_i|y_{1:i-1}, X_{1:n})$

- Search for high probability, tracking top K on each step.
- Beam search is not guaranteed to find optimal solution.
- But much more efficient than exhaustive search.



• Beam search : example (beam size: *K*= 2)



Calculate probability of next word



• Beam search : example (beam size: *K*= 2)



Take top *K* words and compute scores

VestlakeNLP

• Beam search : example (beam size: *K*= 2)



For each of the *K* hypotheses, find top *K* next words and calculate scores

VestlakeNLP

• Beam search : example (beam size: *K*= 2)



Just keep *K* words with highest scores



• Beam search : example (beam size: *K*= 2)

1.9 = P(NLP | < START > he love) + 1.51.5 NLP 0.6 love he study hit • 1.7 = P(study | < START > he love) + 1.11.1 • 2.6 = P(to | < START > I like) + 1.8 $\langle S \rangle$ 1.8 to like play get 0.9 2.7 = P(play | <START > I like) + 1.81.4

For each of the *K* hypotheses, find top *K* next words and calculate scores

VestlakeNLP

• Beam search : example (beam size: *K*= 2)



VestlakeNLP

• Beam search : example (beam size: *K*= 2)



34

VestlakeNLP

• Beam search : example (beam size: *K*= 2)



4.0



• Beam search : example (beam size: *K*= 2)



36
Model 2: LSTM using attention



• Beam search : example (beam size: *K*= 2)



Model 2: LSTM using attention



• Beam search : example (beam size: *K*= 2)



38

Model 2: LSTM using attention



- Beam search: tips
 - Usually we continue beam search until:
 - We reach time step *T* (where *T* is some pre-defined cutoff), or
 - We have at least *n* completed hypotheses (where *n* is pre-defined cutoff)
 - Applying pruning to each time step:
 - Enumeration of all possible candidates results in a *mK*|*V*| time complexity.
 - It is efficient to conduct pruning, keeping only the top 2K candidates.
 - A type value of K in practice is 6.

Contents

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network



- Copying is useful
 - For machine translation, some technological terminologies can be better left.
 - I like to learn NLP. → 我喜欢学习NLP。
 - For dialogue, part of user utterance can be directly included in a system response.
 - Input: *Hello , my name is* **XXX**.
 - Response: *Nice to meet you, XXX*.



- Copying: intuition
 - Basic idea: allow a copying probability to be interpolated with a vocabulary generation probability when generating each word.
 - Source vocabulary *U* (set of words in source sentence) and target vocabulary $V \rightarrow U \cup V$)
 - The integration of copy probabilities into generation probabilities
 - Generation probability:

$$score_g(y_i = \mathbf{v}_j) = \mathbf{p}_i[\mathbf{v}_j], \mathbf{v}_j \in V$$

• How to integrate?



- Copying: How to integrate copy probabilities ?
 - Pointer network:
 - Reuse the attention weights as a distribution over source vocabulary *U*.
 - Copying network:
 - Take a dedicated neural network layer for calculating the copy score.

VestlakeNLP

- Pointer network
 - Take attention score α_{ij} as the distribution of a pointer to a source word:

$$score_c(y_i = x_k) = \alpha_{ik}, k \in [1, ..., n]$$

• The pointers to the same word in the source vocabulary *U* are then aggregated:

$$score_{c}(y_{i} = u_{j}) = \sum_{k:x_{k}=u_{j}} \alpha_{ik}, u_{j} \in U, j \in [1, ..., |U|]$$

• Finally, for the generation of y_i , the generation probability and the copy probability are linearly interpolated:

$$score(y_i) = \lambda score_g(y_i) + (1 - \lambda) score_c(y_i),$$



- Copying network
 - A copying network calculates $score_c$ directly, by taking the current decoding state h_{i-1}^{dec} and the encoder hidden state of a source word h_j^{enc} as input:

$$score_{c}(y_{i} = u_{k}) = \sum_{j:x_{j}=u_{k}} \tanh(\mathbf{h}_{j}^{enc}\mathbf{W}^{c})\mathbf{h}_{i-1}^{dec}$$

• The probability of generating a target word *y_i* is calculated by the sum of a generate probability and a copy probability:

 $P(y_i|X_{1:n}, Y_{1:i-1}) = P_g(y_i|X_{1:n}, Y_{1:i-1}) + P_c(y_i|X_{1:n}, Y_{1:i-1})$

WestlakeNLP

• Copying network

• *P_g* and *P_c* represents normalized versions of *score_g* and *score_c*, which are a part of a distribution *P* of the generation and copy:

$$\begin{split} P_g(y_i = \mathbf{v}_j | X_{1:n}, Y_{1:i-1}) &= \begin{cases} \frac{1}{Z} \exp(score_g(y_i = \mathbf{v}_j)), y_i \in V \\ 0, y_i \notin V \end{cases} \\ P_c(y_i = u_k | X_{1:n}, Y_{1:i-1}) &= \begin{cases} \frac{1}{Z} \exp(score_c(y_i = u_k)), y_i \in U \\ 0, y_i \notin U \end{cases} \end{split}$$

where

$$Z = \sum_{v \in V} \exp(score_g(v)) + \sum_{u \in U} \exp(score_c(u))$$

Contents

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network



- Dealing with OOV words
 - The prefix and the suffix can give us hints about the meaning of a word. Inspired of this, we can make use of subwords to enrich vocabulary instead of a <UNK> tokens in model 1, 2 and 3.
 - First consideration is morphological segmentation, however, it can be ambiguous and not able to cover all words.
 - Byte-pair encoding (BPE) is a useful algorithm:
 Most frequent byte pair → a new byte.



• Original BPE

• Suppose to compress the following text using BPE:

aabaadaab

- The byte-pair *aa* is the most frequent pair, which can be replaced by an unused code *Z*:
 ZbZdZb
- Then repeat the process with replace *Zb* by *Y*:

YZdY

• *YZdY* is the final representation, since there are no pairs that occur more than once.

- Expanding the vocabulary using BPE
 - Start with a vocabulary of characters
 - Most frequent n-gram pairs \rightarrow a new n-gram

Corpus	Vocabulary
2 bed 2 better 3 east 4 west	b, e, d, t, r, a, s, w Start with all characters in vocabulary

- Expanding the vocabulary using BPE
 - Start with a vocabulary of characters
 - Most frequent n-gram pairs \rightarrow a new n-gram



- Expanding the vocabulary using BPE
 - Start with a vocabulary of characters
 - Most frequent n-gram pairs \rightarrow a new n-gram



- Expanding the vocabulary using BPE
 - Start with a vocabulary of characters
 - Most frequent n-gram pairs \rightarrow a new n-gram



- Expanding the vocabulary using BPE
 - Start with a vocabulary of characters
 - Most frequent n-gram pairs \rightarrow a new n-gram

Corpus	Vocabulary	
 2 bed 2 better 3 east 4 west 	b, e, d, t, r, a, s, w, st, be	
	A new OOV word: best \rightarrow be@@ + st	



- Byte Pair Encoding: tips
 - Have a target vocabulary size and stop when you reach it.
 - Segmentation is only within words identified by some prior tokenizer.
 - Automatically decide vocabulary for system.

Contents

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network





WestlakeNLP

• Encoder: Self-attention



Input of a self-attention layer: $\mathbf{V}_{1}^{n} = \mathbf{v}_{1}, \mathbf{v}_{2}, ..., \mathbf{v}_{n}$ $\mathbf{V}_{1}^{n} \in \mathbb{R}^{n \times d_{h}}$ Output of a self-attention layer: $\mathbf{H}_{1}^{n} = \mathbf{h}_{1}, \mathbf{h}_{2}, ..., \mathbf{h}_{n}$ $\mathbf{H}_{1}^{n} \in \mathbb{R}^{n \times d_{h}}$

$$\mathbf{H}_{1}^{n} = attention(\mathbf{V}_{1}^{n}, \mathbf{V}_{1}^{n}, \mathbf{V}_{1}^{n})$$
$$Q = K = V$$

WestlakeNLP

• Encoder: Multi-head self-attention



 $head_{j} = attention(\mathbf{V}_{1}^{n}\mathbf{W}_{j}^{query}, \mathbf{V}_{1}^{n}\mathbf{W}_{j}^{key}, \mathbf{V}_{1}^{n}\mathbf{W}_{j}^{value})$ $\mathbf{h}^{mh} = MultiHead(\mathbf{h}^{in}, \mathbf{h}^{in}, \mathbf{h}^{in})$ $MultiHead(\mathbf{V}_{1}^{n}, \mathbf{V}_{1}^{n}, \mathbf{V}_{1}^{n}) = concat(head_{1}, ..., head_{k})\mathbf{W}^{O}$

WestlakeNLP

• Encoder: Layer normalization



WestlakeNLP

61

• Encoder: Feed-forward sub-layer



WestlakeNLP

• Encoder: output



WestlakeNLP

• Encoder: input encoding and position embedding



WestlakeNLP

• Decoder: similar to encoder



- Obtain the next target word $\mathbf{o}_i = \mathbf{W}^{dec} \mathbf{h}_{i-1}^{dec}$ $\mathbf{p}_i = softmax(\mathbf{o}_i)$
- MLP layer
 - $\mathbf{H}^{dff'} = ReLU(\mathbf{W}^{dff'}\mathbf{H}^{clm} + \mathbf{b}^{dff'})$ $\mathbf{H}^{dff} = \mathbf{W}^{dff}\mathbf{H}^{dff'} + \mathbf{b}^{dff'}$ $\mathbf{H}^{dec} = LayerNorm(\mathbf{H}^{dff} + \mathbf{H}^{clm})$
- Target to source attention
- $\mathbf{H}^{c} = MultiHead(\mathbf{H}^{dlm}, \mathbf{H}^{enc}, \mathbf{H}^{enc})$ $\mathbf{H}^{clm} = LayerNorm(\mathbf{H}^{c} + \mathbf{H}^{dlm})$
- Target self-attention
 - $$\begin{split} \mathbf{H}^{dmh} &= MultiHead(\mathbf{Y}_{1:i-1}, \mathbf{Y}_{1:i-1}, \mathbf{Y}_{1:i-1}) \\ \mathbf{H}^{dlm} &= LayerNorm(\mathbf{Y}_{1:i-1}, \mathbf{H}^{dmh}) \end{split}$$
- Embedding & positional encoding $\mathbf{Y}_{1:i-1} = [emb'(y_1), emb'(y_2), \dots, emb'(y_{i-1})]$

WestlakeNLP

• Training:

Given a set of input-output sentence pairs: $D = \{(X_i, Y_i)\}|_{i=1}^N$

$$X_i = x_1, x_2, \dots, x_{n_i}$$

$$Y_i = y_1^i, y_2^i, ..., y_{m_i}^i$$

Object function:

$$L = -\sum_{i=1}^{N} \sum_{j=1}^{m_i} \log\left(\mathbf{p}_j[y_j^i]\right) = \sum_{i=1}^{N} \sum_{j=1}^{m_i} \log\left(P(y_j^i|X_i, Y_{i:j-1}^i)\right)$$

VestlakeNLP

• Why self attention:

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-attention	$O(n^2 \cdot d)$	0(1)	0(1)
Recurrent	$O(n \cdot d^2)$	O(n)	O(n)
Convolutional	$O(k \cdot n \cdot d^2)$	0(1)	$O(log_k(n))$

n =sequence length d =depth k =kernel size

Contents

VestlakeNLP

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network

• 16.2 Text Matching Models

- 16.2.1 Matching Two Texts in Parallel
- 16.2.2 Searching for a match
- 16.2.3 Memory network

Text Matching Models



• Text matching: whether two pieces of text semantically match each other.

Paraphrase detection, text entailment detection...

• Searching for a match: find a section in a given piece of text that matches the meaning of another piece of text.

Machine reading comprehension...



(a) Matching two pieces of text

(b) Searching for a match

Contents

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network

Text Matching

WestlakeNLP

• Text matching:

How many people live in Melbourne Matching — Matching — Probability

Text Matching

VestlakeNLP

- Siamese network:
 - Apply identical encoders to represent both sentences.
 - Parameters are shared between both encoders.



How many people live in Melbourne What's the population of Melbourne

Text Matching

WestlakeNLP

• Siamese network:



$$\begin{split} &W_{1:n_1}^1 = w_1^1, w_2^1, \dots, w_{n_1}^1 \quad emb(W^1) = [emb(w_1^1); emb(w_2^1); \dots; emb(w_{n_1}^1)] \\ &W_{1:n_2}^2 = w_1^2, w_2^2, \dots, w_{n_2}^2 \quad emb(W^2) = [emb(w_1^2); emb(w_2^2); \dots; emb(w_{n_2}^2)] \end{split}$$
WestlakeNLP

• Siamese network:



$$\begin{split} W_{1:n_1}^1 &= w_1^1, w_2^1, \dots, w_{n_1}^1 \quad emb(W^1) = [emb(w_1^1); emb(w_2^1); \dots; emb(w_{n_1}^1)] \\ W_{1:n_2}^2 &= w_1^2, w_2^2, \dots, w_{n_2}^2 \quad emb(W^2) = [emb(w_1^2); emb(w_2^2); \dots; emb(w_{n_2}^2)] \end{split}$$

VestlakeNLP

• Siamese network:



$$\begin{split} &W_{1:n_1}^1 = w_1^1, w_2^1, \dots, w_{n_1}^1 \quad emb(W^1) = [emb(w_1^1); emb(w_2^1); \dots; emb(w_{n_1}^1)] \\ &W_{1:n_2}^2 = w_1^2, w_2^2, \dots, w_{n_2}^2 \quad emb(W^2) = [emb(w_1^2); emb(w_2^2); \dots; emb(w_{n_2}^2)] \end{split}$$

WestlakeNLP

• Siamese network:



$$\begin{split} &W_{1:n_1}^1 = w_1^1, w_2^1, \dots, w_{n_1}^1 \quad emb(W^1) = [emb(w_1^1); emb(w_2^1); \dots; emb(w_{n_1}^1)] \\ &W_{1:n_2}^2 = w_1^2, w_2^2, \dots, w_{n_2}^2 \quad emb(W^2) = [emb(w_1^2); emb(w_2^2); \dots; emb(w_{n_2}^2)] \end{split}$$

WestlakeNLP

• Training:



$$L = -\sum_{i=1}^{N} (y_i \log P(match(W_i^1, W_i^2)) + (1 - y_i) \log(1 - P(match(W_i^1, W_i^2))))$$

VestlakeNLP

• Attention matching network:



VestlakeNLP

• Attention matching network:



Same to the Siamese network

 $emb(W^{1}) = [emb(w_{1}^{1}); emb(w_{2}^{1}); ...; emb(w_{n_{1}}^{1})]$ $emb(W^{2}) = [emb(w_{1}^{2}); emb(w_{2}^{2}); ...; emb(w_{n_{2}}^{2})]$

VestlakeNLP

• Attention matching network:



$$emb(W^{1}) = [emb(w_{1}^{1}); emb(w_{2}^{1}); ...; emb(w_{n_{1}}^{1})]$$
$$emb(W^{2}) = [emb(w_{1}^{2}); emb(w_{2}^{2}); ...; emb(w_{n_{2}}^{2})]$$

VestlakeNLP

• Attention matching network:



$$s_{ij} = \mathbf{V}^T tanh(\mathbf{W}^1 \mathbf{h}_i^1 + \mathbf{W}^2 \mathbf{h}_j^2 + \mathbf{b})$$
$$\alpha_{ij} = \frac{\exp(s_{ij})}{\sum_{k=1}^{n_2} \exp(s_{ik})}$$
$$\hat{\mathbf{h}}_i^1 = \sum_{j=1}^{n_2} \alpha_{ij}$$

WestlakeNLP

• Attention matching network:



Additive attention Denoted:

$$\widehat{\mathbf{H}}^{1} = sent_att(\mathbf{H}^{1}, \mathbf{H}^{2})$$

$$\widehat{\mathbf{H}}^{2} = sent_att(\mathbf{H}^{2}, \mathbf{H}^{1})$$

WestlakeNLP

• Attention matching network:



Concatenate the original vectors and the matching vectors:

 $\begin{aligned} \mathbf{U} &= avg(\mathbf{Q}^1) \oplus max(\mathbf{Q}^1) \oplus avg(\mathbf{Q}^2) \oplus max(\mathbf{Q}^2) \\ \mathbf{Q}^1 &= [\mathbf{h}_1^1 \oplus \hat{\mathbf{h}}_1^1; \mathbf{h}_2^1 \oplus \hat{\mathbf{h}}_2^1; ...; \mathbf{h}_{n_1}^1 \oplus \hat{\mathbf{h}}_{n_1}^1] \\ \mathbf{Q}^2 &= [\mathbf{h}_1^2 \oplus \hat{\mathbf{h}}_1^2; \mathbf{h}_2^2 \oplus \hat{\mathbf{h}}_2^2; ...; \mathbf{h}_{n_2}^2 \oplus \hat{\mathbf{h}}_{n_2}^2] \end{aligned}$

ave: Bitwise average pooling *max*: Max pooling ⊕: Cancatenate

WestlakeNLP

• Attention matching network:



Final multilayer perceptron for matching probability: $P(match(W^1, W^2)) = \sigma(MLP(\mathbf{U}))$

 $\mathbf{U} = avg(\mathbf{Q}^{1}) \oplus max(\mathbf{Q}^{1}) \oplus avg(\mathbf{Q}^{2}) \oplus max(\mathbf{Q}^{2})$ $\mathbf{Q}^{1} = [\mathbf{h}_{1}^{1} \oplus \hat{\mathbf{h}}_{1}^{1}; \mathbf{h}_{2}^{1} \oplus \hat{\mathbf{h}}_{2}^{1}; ...; \mathbf{h}_{n_{1}}^{1} \oplus \hat{\mathbf{h}}_{n_{1}}^{1}]$ $\mathbf{Q}^{2} = [\mathbf{h}_{1}^{2} \oplus \hat{\mathbf{h}}_{1}^{2}; \mathbf{h}_{2}^{2} \oplus \hat{\mathbf{h}}_{2}^{2}; ...; \mathbf{h}_{n_{2}}^{2} \oplus \hat{\mathbf{h}}_{n_{2}}^{2}]$

ave: Bitwise average pooling max: Max pooling ⊕: Cancatenate

VestlakeNLP

- Bidirectional attention matching network:
 - Similarly, Bi-direction attention (*co-attention*) can be used to extract the mutual information between H^1 and H^2 .

$$\widetilde{\mathbf{H}}^{1} = Dup(\mathbf{H}^{1}, n_{2})$$
$$\widetilde{\mathbf{H}}^{2} = Dup(\mathbf{H}^{2}, n_{1})$$
$$\mathbf{S} = avg_{1}(\mathbf{V}(\widetilde{\mathbf{H}}^{1} \bigoplus (\widetilde{\mathbf{H}}^{2})T_{2,3} \bigoplus (\widetilde{\mathbf{H}}^{1} \otimes (\widetilde{\mathbf{H}}^{2})T_{2,3})))$$

• Attention weight matrix:

 $\alpha_1 = softmax_2(\mathbf{S}) \\ \alpha_2 = softmax_2(\mathbf{S}^{\mathsf{T}})$

• The bi-directional attention output **U**

$$\mathbf{U} = [\mathbf{H}^1 \oplus \mathbf{H}^2 \oplus [\mathbf{H}^1 \otimes oldsymbol{lpha}_1] \oplus [\mathbf{H}^2 \otimes oldsymbol{lpha}_2]]$$

Contents

VestlakeNLP

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network

VestlakeNLP

• Network structure searching for a match:



87

Searching for a Match

START

Given a context sequence W^1 and a query W^2 , the task is to find an answer span

$$Y = (W^1)_b^e = w_b^1, \dots, w_e^1$$

that matches W^2 .

• Network structure searching for a match:



END



• Network structure searching for a match:



Encoding with Bi-LSTM

 $\mathbf{H}^{1} = BiLSTM(emb(W^{1}))$ $\mathbf{H}^{2} = BiLSTM(emb(W^{2}))$



• Network structure searching for a match:



Additive attention for matching:

$$s_{ij} = \mathbf{V}^T tanh(\mathbf{W}^1 \mathbf{h}_i^1 + \mathbf{W}^2 \mathbf{h}_j^2 + \mathbf{b})$$
$$\boldsymbol{\alpha}_{ij} = \frac{\exp(s_{ij})}{\sum_{k=1}^{n_2} \exp(s_{ik})}$$
$$\hat{\mathbf{h}}_i^1 = \sum_{j=1}^{n_2} \alpha_{ij} \mathbf{h}_j^2$$

• Network structure searching for a match:



Final representation:

 $\mathbf{U}^1 = [\mathbf{h}_1^1 \oplus \hat{\mathbf{h}}_1^1; \mathbf{h}_2^1 \oplus \hat{\mathbf{h}}_2^1; \dots; \mathbf{h}_{n_1}^1 \oplus \hat{\mathbf{h}}_{n_1}^1]$

WestlakeNLP

• Network structure searching for a match:



Start Predicting:

$$\mathbf{M}^{b} = BiLSTM(\mathbf{U}^{1})$$
$$\mathbf{Q}^{b} = [\mathbf{u}_{1}^{1} \bigoplus \mathbf{m}_{1}^{b}; \mathbf{u}_{2}^{1} \bigoplus \mathbf{m}_{2}^{b}; ...; \mathbf{u}_{n_{1}}^{1} \bigoplus \mathbf{m}_{n_{1}}^{b}]$$
$$\mathbf{s}^{b}[i] = (\mathbf{v}^{b})^{T}\mathbf{Q}_{i}^{b}, i \in [1, ..., n_{1}]$$
$$\mathbf{p}^{b} = softmax(\mathbf{s}^{b})$$



92

Searching for a Match

• Network structure searching for a match:



End Predicting:

$$\mathbf{Q}^{e} = [\mathbf{u}_{1}^{1} \bigoplus \mathbf{m}_{1}^{e}; \mathbf{u}_{2}^{1} \bigoplus \mathbf{m}_{2}^{e}; ...; \mathbf{u}_{n_{1}}^{1} \bigoplus \mathbf{m}_{n_{1}}^{e}]$$

$$s^{e}[i] = (\mathbf{v}^{e})^{T} \mathbf{q}_{i}^{e}, i \in [1, ..., n_{1}]$$

$$\mathbf{p}^{e} = softmax(\mathbf{s}^{e})$$

WestlakeNLP

• Network structure searching for a match:



Training:

$$L = -\frac{1}{N} \sum_{i}^{N} (\log(\mathbf{p}^{b}[b_{i}]) + \log(\mathbf{p}^{e}[e_{i}]))$$

Contents

VestlakeNLP

- 16.1 Sequence to sequence models
 - 16.1.1 Model 1: Seq2seq using LSTM
 - 16.1.2 Model 2: Adding more source features using target-tosource attention
 - 16.1.3 Model 3: Copying form the source
 - 16.1.4 Subword Embedding
 - 16.1.5 Seq2seq using Attention Network
- 16.2 Text Matching Models
 - 16.2.1 Matching Two Texts in Parallel
 - 16.2.2 Searching for a match
 - 16.2.3 Memory network



- Previous methods are weak in some case:
 - In many case, the correct answer to a question cannot be directly read off one evidence, but multi-stage inference is required instead.
 - For instance, suppose that the document is

Joe travelled to the office. Joe left the milk. Joe went to the bathroom. Q: Where is the milk now?

one needs to find out the most relevant fact: *"Joe left the milk"* Then, one needs to consider the second relevant fact, *"Joe travelled to the office"*, before finding a to describe the answer *"office"*.



- Why memory network:
 - Neural networks have hard times to capture long-range dependencies, even LSTMs.
 - Multi-stage inference is necessary in machine reading comprehension.
 - Memory networks try to overcome these problem using an external memory.



- Intuition of memory network:
 - A neural network: LSTM, SAN, ...



VestlakeNLP

- Intuition of memory network:
 - A neural network: LSTM, SAN, ...
 - An external memory





WestlakeNLP

- Intuition of memory network:
 - A neural network: LSTM, SAN, ...
 - An external memory where the neural network can write and read to





- Memory network structure: a memory and 4 components
 - Input feature map(I): encodes the input into a hidden representations, e.g. bag of words
 - Generalization(G): stores the input into memory
 - Output feature map(O): read the most relevant memory slots
 - Response(R): given the info. read from the memory, return the output.



WestlakeNLP

â

• End2end memory network for QA: an example

Joe travelled to the office. Joe left the milk. Joe went to the bathroom. Q: Where is the milk now?





Component I: Input Representation

Raw text sentence is transformed in its vector representation e.g. BOW.

$$\mathbf{m}_{i} = \sum_{j=1}^{n_{i}} emb^{\mathbf{E}^{R}}\left(w_{j}^{i}\right)$$



Joe travelled to the office. Joe left the milk. Joe went to the bathroom.

 l_j is position encoding to capture the order of the words. $l_j^i = (1 - j/n_i) - (i/d_h)(1 - 2j/n_i)$

VestlakeNLP

- Component G:
 - The sentences are then written to the memory sequentially, via the component G:



• Notice that the memory is fixed in this approach once is written, it is not changed neither during learning nor during testing.

VestlakeNLP

- Query text Representation
 - Similar, query representation as follow:

$$\mathbf{u} = \sum_{i=1}^{n_Q} l_i^Q \cdot emb^{\mathbf{E}^Q}(W_i^Q) a \qquad \bigcirc \bullet \bullet \bullet \bullet \qquad Q: Where is the milk now?$$

 l_j is position encoding to capture the order of the words.

$$l_i^Q = (1 - i/n_Q) - (1/d_h)(1 - 2i/n_Q)$$

VestlakeNLP

- Component O: matching
 - The most relevant memory blocks are found by calculating p over memory cells m_i :



VestlakeNLP

• Component O: output

• Then, we calculate a context vector *o* in accordance with *p*.





- Component R: response with component O
 - *a*[*i*] represents the probability of the *i* th word from the vocabulary being the answer word with the evidence vector **o**.



 $\mathbf{q} = \mathbf{W}(\mathbf{o} + \mathbf{u})$ $\hat{\mathbf{a}} = softmax(\mathbf{q})$

• But, we can not get the answer directly. More steps can be necessary.





- Component R: response with multi-hop
 - We can leverage the existing evidence vector o for finding a new distribution p over memory cells, which in turn results in a new context representation of o.



 $\mathbf{u}_{k+1} = \mathbf{u}_k + \mathbf{o}_k$ $\mathbf{q}_K = \mathbf{W}(\mathbf{o}_K + \mathbf{u}_K)$

• This is typically called a *hop* in inference.
Memory Network

WestlakeNLP

• Training

• Given a set of training data: $D = \{(D_i^R, W_i^Q, a_i)\}|_{i=1}^N$



Standard cross-entropy loss with the inference above:

$$L = -\sum_{i}^{N} \log\left(\hat{\mathbf{a}}_{i}[a_{i}]\right)$$