NLP1 2020 - Lecture 7

Compositional semantics and sentence representations

Mario Giulianelli

Institute for Logic, Language and Computation



18 November 2020

Outline

Compositional semantics

Compositional distributional semantics

Compositional semantics with neural networks

Outline

Compositional semantics

Compositional distributional semantics

Compositional semantics with neural networks

3

Principle of compositionality

Principle of compositionality



Principle of compositionality





Principle of compositionality



Semantic composition

What is the meaning of "carnivorous plants digest slowly"?



- Similar syntactic structures may have different meanings
 - ➡ it runs
 - → it rains, it snows (here, it is a pleonastic pronoun)
- Different syntactic structures may have the same meaning (e.g., passive constructions)
 - ➡ Eve ate the apple.
 - ➡ The apple was eaten by Eve.
- Not all phrases are interpreted compositionally (e.g., idioms)
 - kick the bucket
 - ➡ pull someone's leg

but the compositional interpretation is still possible.

- Additional meaning can arise through composition (e.g., logical metonymy)
 - ➡ fast car
 - → fast algorithm
 - ➡ begin a book
- Meaning transfers (e.g., metaphor)
 - → he put a grape into his mouth and <u>swallowed</u> it whole
 - → he <u>swallowed</u> her story whole
- Additional connotations can arise through composition
 - I can't <u>buy</u> this story
- Recursive composition

Recursive composition



"Of course I care about how you imagined I thought you perceived I wanted you to feel."

Recursive composition



"Of course I care about how you imagined I thought you perceived I wanted you to feel."

A defining property of natural languages is **productivity**: they license a theoretically infinite set of possible expressions.

Compositionality and recursion allow for productivity.

Cautionary notes

- The meaning of the whole is constructed from its parts, and the meaning of the parts is derived from the whole.
- Compositionality is a matter of degree rather than a binary notion.



Modelling compositional semantics

1. Compositional distributional semantics

- composition is modelled in a vector space
- ► unsupervised
- general purpose representations

2. Compositional semantics with **neural networks**

- ► (typically) supervised
- (typically) task-specific representations

Outline

Compositional semantics

Compositional distributional semantics

Compositional semantics with neural networks

it was authentic <u>scrumpy</u>, rather sharp and very strong we could taste a famous local product — <u>scrumpy</u> spending hours in the pub drinking <u>scrumpy</u>

it was authentic <u>scrumpy</u>, rather sharp and very strong we could taste a famous local product — <u>scrumpy</u> <i>spending hours in the pub drinking <u>scrumpy</u>

Can distributional semantics can be extended to account for the meaning of phrases and sentences?

it was authentic <u>scrumpy</u>, rather sharp and very strong we could taste a famous local product — <u>scrumpy</u> spending hours in the pub drinking <u>scrumpy</u>

Can distributional semantics can be extended to account for the meaning of phrases and sentences?

her <u>old dog</u> is turning 14 this year! you see your <u>old dog</u> lumber slowly to the food bowl in an <u>old dog</u>, behaviour changes can appear suddenly

Can distributional semantics can be extended to account for the meaning of phrases and sentences?

- Given a finite vocabulary, natural languages licence an infinite amount of sentences.
- So it is impossible to learn vector representations for all sentences.

Can distributional semantics can be extended to account for the meaning of phrases and sentences?

- Given a finite vocabulary, natural languages licence an infinite amount of sentences.
- So it is impossible to learn vector representations for all sentences.
- But we can still use distributional word representations and learn to perform semantic composition in distributional space.

Can distributional semantics can be extended to account for the meaning of phrases and sentences?

Given a finite vocabulary, natural languages licence an infinite amount of sentences.

21

- So it is impossible to learn vector representations for all sentences.
- But we can still use distributional word representations and learn to perform semantic composition in distributional space.

vector mixture models lexical function models

Mitchell and Lapata. (2010). Composition in distributional models of semantics. Cognitive science.

Words are represented as vectors, which combine to produce new vectors.

p = f(u, v, R, K)

Mitchell and Lapata. (2010). Composition in distributional models of semantics. Cognitive science.

Words are represented as vectors, which combine to produce new vectors.

p = f(u, v, R, K)

Mitchell and Lapata. (2010). Composition in distributional models of semantics. Cognitive science.

Words are represented as vectors, which combine to produce new vectors.

p = f(u, v, R, K)

Mitchell and Lapata. (2010). Composition in distributional models of semantics. Cognitive science.

Words are represented as vectors, which combine to produce new vectors.

p = f(u, v, R, K)

<u>Constraint</u>: **p** lies in the same n-dimensional space as **u** and **v**. <u>Assumption</u>: all syntactic types are similar enough to have the same dimensionality.

Mitchell and Lapata. (2010). Composition in distributional models of semantics. Cognitive science.

Words are represented as vectors, which combine to produce new vectors.

p = f(u, v, R, K)

<u>Constraint</u>: **p** lies in the same n-dimensional space as **u** and **v**. <u>Assumption</u>: all syntactic types are similar enough to have the same dimensionality.



NLP1 2020: Compositional semantics and sentence representations

Additive and multiplicative models



				addi	tive	multiplicative	
	\mathbf{dog}	\mathbf{cat}	old	$\mathbf{old} + \mathbf{dog}$	$\mathbf{old} + \mathbf{cat}$	$\mathbf{old} \odot \mathbf{dog}$	$\mathbf{old} \odot \mathbf{cat}$
runs	1	4	0	1	4	0	0
barks	5	0	7	12	7	35	0

Additive and multiplicative models

- The additive and the multiplicative model are symmetric (commutative): they do not take word order or syntax into account.
 - → John hit the ball = The ball hit John
- Correlate with human similarity judgments about adjective-noun, noun-noun, verb-noun and noun-verb pairs
- More suitable for modelling content words, would not apply well to function words:
 - ➡ <u>some</u> dogs, lice <u>and</u> dogs, lice <u>on</u> dogs

Lexical function models

Assumption: all syntactic types are similar enough to have the same dimensionality.

p = f(U, v, R, K)

Lexical function models

Assumption: all syntactic types are similar enough to have the same dimensionality.

 $\mathbf{p} = f(\mathbf{U}, \mathbf{v}, R, K)$

Distinguish between:

- words whose meaning is directly determined by their distributional profile, e.g. nouns
- words that act as functions transforming the distributional profile of other words, e.g., adjectives, adverbs

 $\mathbf{p} = f(\mathbf{U}, \mathbf{v}, ADJ) = \mathbf{U}\mathbf{v}$



Lexical function models

Baroni and Zamparelli. (2010). Nouns are vectors, adjectives are matrices: Representing adjectivenoun constructions in semantic space. In *Proceedings of EMNLP*.

Adjectives modelled as lexical functions that are applied to nouns: *old dog = old(dog)*

- ► Adjectives are parameter matrices (Aold, Afurry, etc.)
- ► Nouns are vectors (house, dog, etc.)
- Composition is a linear transformation: **old dog** = $A_0/d \times dog$.

Learning adjective matrices

For each adjective, learn a parameter matrix that allows to predict adjective-noun phrase vectors.

	Х	Y
Training set	house dog car cat toy	old house old dog old car old cat old toy
Test set	elephant mercedes	old elephant old mercedes

Learning adjective matrices

- 1. Obtain a distributional vector \mathbf{n}_j for each noun n_j in the vocabulary using a conventional DSM.
- 2. Collective all adjective-noun pairs (a_i, n_j) from the corpus.
- 3. Obtain a distributional phrase vector \mathbf{p}_{ij} for each pair (a_i, n_j) from the same corpus using a conventional DSM—treating the phrase $a_i n_j$ as a single word.
- 4. The set of tuples $\{(\mathbf{n}_i, \mathbf{p}_{ij})\}_i$ represents a dataset $\mathcal{D}(a_i)$ for the adjective a_i .
- 5. Learn matrix A_i from $\mathcal{D}(a_i)$ using linear regression. Minimise the squared error loss:

$$L(\mathbf{A}_i) = \sum_{j \in \mathcal{D}(a_i)} \|\mathbf{p}_{ij} - \mathbf{A}_i \mathbf{n}_j\|^2$$

Verbs as lexical functions

Verbs too can modelled as lexical functions that are applied to their arguments.

They are represented as tensors whose order is determined by the subcategorisation frame of the verb (i.e., how many and what type of arguments the verb takes).

Intransitive verbs take a subject as their only argument

dogs <u>bark</u> Vbark \times dogs

modelled as a matrix (second-order tensor)

Transitive verbs take a subject and an object

dogs eat meat $(V_{eat} \times meat) \times dogs$

modelled as a third-order tensor

Modelling compositional semantics

1. Compositional distributional semantics

- composition is modelled in a vector space
- unsupervised learning
- general purpose representations

2. Compositional semantics with neural networks

- (typically) supervised learning
- (typically) task-specific representations

Modelling compositional semantics

1. Compositional distributional semantics

- composition is modelled in a vector space
- unsupervised learning
- general purpose representations

2. Compositional semantics with neural networks

- (typically) supervised learning
- (typically) task-specific representations

Task: sentiment classification (Practical 2)

36
Outline

Compositional semantics

Compositional distributional semantics

Compositional semantics with neural networks

Compositional semantics with NNs

- 1. Learn sentence (or phrase) representations
- 2. Learn to make task-specific predictions based on the sentence (or phrase) representation

Compositional semantics with NNs

- 1. Learn sentence (or phrase) representations
- 2. Learn to make task-specific predictions based on the sentence (or phrase) representation

Darkly funny and frequently insightful

0. very negative

1. negative

2. neutral

3. positive

Darkly funny and frequently insightful

0. very negative

1. negative

2. neutral

3. positive

Darkly funny and frequently insightful

0. very negative

1. negative

2. neutral

3. positive

Darkly funny and frequently insightful

0. very negative

1. negative

2. neutral

3. positive

Darkly funny and frequently insightful

0. very negative

- 1. negative
- 2. neutral
- **3. positive**
- 4. very positive

Darkly funny and frequently insightful

- 0. very negative
- 1. negative
- 2. neutral
- **3. positive**
- 4. very positive

Task-specific representations



Dataset: Stanford Sentiment Treebank

13

12K movie reviews

16

Darkly

- one sentence per review
- sentence-level sentiment score
- binary syntactic tree
- phrase-level sentiment scores

13

and

funny

13

frequently

151/3

insightful

Compositional semantics with NNs

1. Learn sentence (or phrase) representations

2. Learn to make task-specific predictions based on the sentence (or phrase) representation

Compositional semantics with NNs

Models

- 1. Bag of Words (BOW)
- 2. Continuous Bag of Words (CBOW)
- 3. Deep Continuous Bag of Words (Deep CBOW)
- 4. Deep CBOW with pre-trained word embeddings
- 5. LSTM
- 6. Tree-LSTM

- Additive model: does not take word order or syntax into account
- Task-specific word representations with **fixed dimensionality** (d = 5)

50

Dimensions of vector space are explicit, interpretable

Sum word embeddings, add bias



Sum word embeddings, add bias



Sum word embeddings, add bias



this	[0.0,	0.1,	0.1,	0.1,	0.0]
movie	[0.0,	0.1,	0.1,	0.2,	0.1]
is	[0.0,	0.1,	0.0,	0.0,	0.0]
stupid	[0.9,	0.5,	0.1,	0.0,	0.0]

bias	[0.0,	0.0,	0.0,	0.0,	0.0]
sum	[0.9,	0.8,	0.3,	0.3,	0.1]

argmax: 0 (very negative)

- Additive model: does not take word order or syntax into account
- Task-specific word representations of arbitrary dimensionality
- Dimensions of vector space are not interpretable
- Prediction can be traced back to the sentence vector dimensions

Sum word embeddings, project to 5D using W, add bias: W ($\sum x_{t}$) + **b**







Variable sentence vector size, dependent on sentence length

- Not very sensible conceptually
 - sentences in a different vector space than words
 - one vector space for each sentence length in the dataset
- Hardly practicable
 - ➡ what size should the transformation matrix be?
 - vector size can grow very large

- Additive model: does not take word order or syntax into account
- Task-specific word representations of arbitrary dimensionality
- Dimensions of vector space are not interpretable
- More layers and non-linear transformations: prediction cannot be easily traced back

W" tanh(W ($\sum \mathbf{x}_t$) + **b**) + **b**') + **b**"



Deep CBOW + pre-trained embeddings

- Additive model: does not take word order or syntax into account
- Pre-trained general-purpose word representations (e.g., Skip-gram, GloVe)
 - → frozen: not updated during training
 - → fine-tuned: updated with task-specific learning signal (specialised)
- Dimensions of vector space are not interpretable
- Multiple layers and non-linear transformations: prediction cannot be easily traced back

Deep CBOW + pre-trained embeddings

- Additive model: does not take word order or syntax into account
- Pre-trained general-purpose word representations (e.g., Skip-gram, GloVe)
 - → frozen: not updated during training
 - ➡ fine-tuned: updated with task-specific learning signal (specialised)
- Dimensions of vector space are not interpretable
- Multiple layers and non-linear transformations: prediction cannot be easily traced back



Deep CBOW + pre-trained embeddings

W" tanh(W ($\sum \mathbf{x}_t$) + **b**) + **b**') + **b**"



Training feedforward models

Train networks using Stochastic Gradient Descent (SGD):

- 1. Sample a training example
- 2. Forward pass: compute network activations and obtain an output vector
- 3. Loss: compare output vector with ground-truth label
- 4. Compute gradient of loss w.r.t (trainable) network parameters
- 5. **Backpropagation**: update parameters taking a step in the opposite direction of the gradient

output vector (logits)

o = [-0.1, 0.1, 0.1, 2.4, 0.2]





 $\hat{y} = [0.06, \, 0.07, \, 0.07, \, 0.72, \, 0.08]$



 $\hat{y} = [0.06, 0.07, 0.07, 0.72, 0.08]$

$$y = [0, 0, 0, 1, 0]$$



NLP1 2020: Compositional semantics and sentence representations

Recurrent neural networks

NLP1 2020: Compositional semantics and sentence representations

Recurrent neural networks

- Networks that contain a cycle within their connections
 - The value of a network unit is dependent on earlier outputs
- Naturally handle sequential input: process one element at a time
- Drop history independence assumption
 - Capture and exploit the temporal nature of language
 - Capture word order (and syntactic structure)

Recurrent neural networks



72
Elman, J. L. (1990). Finding structure in time. Cognitive science.



Elman, J. L. (1990). Finding structure in time. Cognitive science.



Elman, J. L. (1990). Finding structure in time. Cognitive science.



Darkly funny and frequently insightful x_1 x_2 x_3 x_4 x_5

$$h_1 = f(x_1, h_0)$$

$$h_2 = f(x_2, f(x_1, h_0)) = f(x_2, h_1)$$

$$h_3 = f(x_3, f(x_2, f(x_1, h_0))) = f(x_3, h_2)$$

Elman, J. L. (1990). Finding structure in time. Cognitive science.



77

Elman, J. L. (1990). Finding structure in time. Cognitive science.

- Consist of a single cell, the "RNN cell", that is fed input elements one at a time (one for each timestep)
- Hidden layer from previous timestep provides a form of memory of preceding context
 - encodes earlier processing
 - informs the decisions to be made at later points in time
- No fixed limit on the length of the preceding context



Elman, J. L. (1990). Finding structure in time. Cognitive science.



Ŵ

 $h_t = f(Uh_{t-1} + Wx_t)$

Elman, J. L. (1990). Finding structure in time. Cognitive science.



$$h_t = f(Uh_{t-1} + Wx_t)$$

$$y_t = g(Vh_t)$$

Elman, J. L. (1990). Finding structure in time. Cognitive science.

$$x \leftarrow (x_1, \dots, x_n)$$

$$h_0 \leftarrow 0$$

for $i \leftarrow 1$ to n do

$$h_t = f(Uh_{t-1} + Wx_t)$$

$$y_t = g(Vh_t)$$

return y



Elman, J. L. (1990). Finding structure in time. Cognitive science.



timesteps

Elman, J. L. (1990). Finding structure in time. Cognitive science.



timesteps

Simple RNN: vanishing gradient

Simple RNNs are hard to train because of the vanishing gradient problem:

- during backpropagation, gradients can quickly become small
- as they repeatedly go through multiplications and non-linear functions (e.g. sigmoid or tanh)



83

LSTM: long short-term memory

Hochreiter and Schmidhuber (1997). Long short-term memory. Neural Computation.

- Designed to maintain relevant context over time
 - learn to forget information that is no longer needed
 - Iearn to remember information required for future decisions
- Deals well with long-term dependencies in the input sequence
- With respect to a Simple RNN
 - → add an explicit context layer (*cell state*) to the architecture
 - → make use of gates to control the flow of information
- Consists of a single cell, the "LSTM cell", that is repeatedly applied to the input elements

LSTM: long short-term memory

Hochreiter and Schmidhuber (1997). Long short-term memory. Neural Computation.



85

LSTM: forget gate

Delete information from the cell state that is no longer needed.

 $f_t = \sigma(U^f h_{t-1} + W^f x_t)$



86

LSTM: forget gate

Delete information from the cell state that is no longer needed.

 $f_t = \sigma(U^f h_{t-1} + W^f x_t)$



LSTM: candidate cell

Extract information from the previous hidden state and the current input.

$$\tilde{c}_t = \tanh(U^c h_{t-1} + W^c x_t)$$



LSTM: input gate

Select the information to add to the new cell state.

 $i_t = \sigma(U^i h_{t-1} + W^i x_t)$



LSTM: input gate

Select the information to add to the new cell state.

 $i_t = \sigma(U^i h_{t-1} + W^i x_t)$



LSTM: cell state

Update the cell state.





LSTM: output gate

Select what information is required for the current hidden state (as opposed to what information needs to be preserved for future decisions).

$$o_t = \sigma(U^o h_{t-1} + W^o x_t)$$



LSTM: hidden state

Update the hidden state.

 $h_t = o_t \odot \tanh(c_t)$



LSTM: all equations

$$f_{t} = \sigma(U^{f}h_{t-1} + W^{f}x_{t})$$

$$i_{t} = \sigma(U^{i}h_{t-1} + W^{i}x_{t})$$

$$o_{t} = \sigma(U^{o}h_{t-1} + W^{o}x_{t})$$

$$\tilde{c}_{t} = \tanh(U^{c}h_{t-1} + W^{c}x_{t})$$

$$c_{t} = f_{t} \odot \tilde{c}_{t-1} + i_{t} \odot \tilde{c}_{t}$$

$$h_{t} = o_{t} \odot \tanh(c_{t})$$



LSTM: predictions



LSTM: stacked and bidirectional



96

LSTM: applications

- Language modelling (Mikolov et al., 2010; Sundermeyer et al., 2012)
- Sequence labelling (e.g., POS tagging, Named Entity Recognition)
- **Parsing** (Vinyals et al., 2015; Kiperwasser and Goldberg, 2016; Dyer et al., 2016)
- Machine translation (Bahdanau et al., 2015)
- ► Auto-regressive generation (e.g., image captioning (Bernardi et al., 2016))
- Sequence classification (e.g., sentiment analysis, fake news detection)

▶ ...

Sentence representations with NNs

Bag of Words models

sentence representations are order-independent function of the word representations

Sequence models

- sentence representations are an order-sensitive function of a sequence of word representations (surface form)
- can they still capture the underlying syntactic structure? (e.g., Linzen et al., 2016, Giulianelli et al., 2018)

Tree-structured models

sentence representations are a function of the word representations, sensitive to the syntactic structure of the sentence

Tree-structured models

- More faithful operationalisation of the principle of compositionality
- Helpful in disambiguation: similar surface form, different underlying structure
- Requires syntactic parse trees (constituency or dependency)



Tree-LSTM

A generalisation of the LSTM to tree-structured input.

- Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved Semantic Representations From Tree-Structured Long Short-Term Memory Networks. ACL 2015.
 - ➡ Child-Sum Tree-LSTM
 - ➡ N-ary Tree-LSTM
- Phong Le and Willem Zuidema. Compositional distributional semantics with long short term memory. *SEM 2015.
- Xiaodan Zhu, Parinaz Sobihani, and Hongyu Guo. Long short-term memory over recursive structures. ICML 2015.

Tree-LSTM



Tree-LSTM

- Gates and memory cell updates are dependent on the states of a node's children, rather than on the states of the previous words.
- Instead of a single forget gate, Tree-LSTM unit contains one forget gate for each child to selectively incorporate information from each child.



Tree-LSTM variants

Child-Sum Tree-LSTM

- → sum over the hidden representations of all children of a node (**no children order**)
- can be used for a variable number of children
- shares parameters between children
- suitable for dependency trees

► N-ary Tree-LSTM

- discriminates between children node positions (weighted sum)
- fixed maximum branching factor: can be used with N children at most
- different parameters for each child
- suitable for constituency trees

Child-Sum Tree-LSTM

$$\tilde{h}_{j} = \sum_{k \in C(j)} h_{k}$$

$$f_{jk} = \sigma(U^{f}h_{k} + W^{f}x_{j})$$

$$i_{j} = \sigma(U^{i}\tilde{h}_{j} + W^{i}x_{j})$$

$$o_{j} = \sigma(U^{o}\tilde{h}_{j} + W^{o}x_{j})$$

$$\tilde{c}_{j} = \tanh(U^{c}\tilde{h}_{j} + W^{c}x_{j})$$

$$c_{j} = i_{j} \odot \tilde{c}_{j} + \sum_{k \in C(j)} f_{jk} \odot c_{k}$$

$$h_{j} = o_{j} \odot \tanh(c_{j})$$



Child-Sum Tree-LSTM

$$\tilde{h}_{j} = \sum_{k \in C(j)} h_{k}$$

$$f_{jk} = \sigma(U^{f}h_{k} + W^{f}x_{j})$$

$$i_{j} = \sigma(U^{i}\tilde{h}_{j} + W^{i}x_{j})$$

$$o_{j} = \sigma(U^{o}\tilde{h}_{j} + W^{o}x_{j})$$

$$\tilde{c}_{j} = \tanh(U^{c}\tilde{h}_{j} + W^{c}x_{j})$$

$$c_{j} = i_{j} \odot \tilde{c}_{j} + \sum_{k \in C(j)} f_{jk} \odot c_{k}$$

$$h_{j} = o_{j} \odot \tanh(c_{j})$$



N-ary Tree-LSTM

$$\begin{split} f_{jk} &= \sigma \left(W^{f} x_{j} + \sum_{l=1}^{N} U_{kl}^{f} h_{jl} \right) \\ i_{j} &= \sigma \left(W^{i} x_{j} + \sum_{l=1}^{N} U_{l}^{i} h_{jl} \right) \\ o_{j} &= \sigma \left(W^{o} x_{j} + \sum_{l=1}^{N} U_{l}^{o} h_{jl} \right) \\ \tilde{c}_{j} &= \tanh \left(W^{c} x_{j} + \sum_{l=1}^{N} U_{l}^{c} h_{jl} \right) \\ \tilde{c}_{j} &= i_{j} \odot \tilde{c}_{j} + \sum_{l=1}^{N} f_{jl} \odot c_{l} \end{split}$$

N-ary Tree-LSTM

$$f_{jk} = \sigma \left(W^f x_j + \sum_{l=1}^N U^f_{kl} h_{jl} \right)$$
$$i_j = \sigma \left(W^i x_j + \sum_{l=1}^N U^i_l h_{jl} \right)$$
$$o_j = \sigma \left(W^o x_j + \sum_{l=1}^N U^o_l h_{jl} \right)$$
$$\tilde{c}_j = \tanh \left(W^c x_j + \sum_{l=1}^N U^c_l h_{jl} \right)$$
$$c_j = i_j \odot \tilde{c}_j + \sum_{l=1}^N f_{jl} \odot c_l$$



 $4 + 3N + N^2$ parameter matrices

Tree-LSTM: forget gates

Child-Sum Tree-LSTM

 $f_{j,left} = \sigma(U^{f}h_{left} + W^{f}x_{j})$ $f_{j,right} = \sigma(U^{f}h_{right} + W^{f}x_{j})$

1 parameter matrix U^f : symmetric treatment of children

Binary Tree-LSTM

$$f_{j,left} = \sigma \left(W^f x_j + U^f_{left,left} h_{j,left} + U^f_{left,right} h_{j,right} \right)$$

$$f_{j,right} = \sigma \left(W^f x_j + U^f_{right,left} h_{j,left} + U^f_{right,right} h_{j,right} \right)$$

4 parameter matrices $U_{x,y}^{f}$: asymmetric treatment of children
Binary Tree-LSTM: sentiment classification



Binary Tree-LSTM: sentiment classification



Compositional semantics with NNs

Models

- 1. Bag of Words (BOW)
- 2. Continuous Bag of Words (CBOW)
- 3. Deep Continuous Bag of Words (Deep CBOW)
- 4. Deep CBOW with pre-trained word embeddings
- 5. LSTM
- 6. Tree-LSTM