

Microsoft Research 微软亚洲研究院

Adaptive Multi-Compositionality for Recursive Neural Models with Applications to Sentiment Analysis

Li Dong^{†*} Furu Wei[‡] Ming Zhou[‡] Ke Xu[†]

[†]State Key Lab of Software Development Environment, Beihang University, Beijing, China [‡]Microsoft Research, Beijing, China donglixp@gmail.com {fuwei,mingzhou}@microsoft.com kexu@nlsde.buaa.edu.cn



July 31, 2014

Semantic Composition

- Principle of Compositionality
 - The meaning of a complex expression is determined by the meanings of its constituent expressions and the rules used to combine them
- Compositional nature of natural language
 - Go beyond words towards sentences

Examples

. . .

- red car -> red + car
- not very good -> not + (very + good)
- eat food -> eat + food

Recursive Neural Models (RNMs)

- Utilize the recursive structures of sentences to obtain the semantic representations
 - The vector representations are used as features and fed into a softmax classifier to predict their labels
- Learn to recursively perform semantic compositions in vector space
- One family of the popular deep learning models



Semantic Composition with Matrix/Tensor

 The main difference among the recursive neural models (RNMs) lies in semantic composition methods

$$v = f\left(\begin{bmatrix} \bullet + \bullet + \bullet + \bullet \\ \bullet + \bullet + \bullet + \bullet \end{bmatrix} \right)$$

$$\boldsymbol{v} = f\left(\boldsymbol{W}\begin{bmatrix}\boldsymbol{v}_l\\\boldsymbol{v}_r\end{bmatrix} + \boldsymbol{b}\right)$$

RNN (Socher et al. 2011)



Problem: RNN and RNTN employ the same global composition function for all pair of input vectors

Motivation of This Work

- Use different composition functions for different types of compositions
 - Negation: not good, not bad

•••

- Intensification: very good, pretty bad
- Contrast: the movie is good, but I love it
- Sentiment word + target/aspect: good movie, low price
- Model the composition as a distribution over multiple composition functions, and adaptively select them



Figure 1: Composition process for "*not so good*". The *g* is a global composition function in recursive neural models.

One Global Composition Function



Figure 2: The composition pool consists multiple composition functions. It selects the functions depending on the input child vectors, and produces the composition result using more than one composition functions.

Adaptive Multi-Compositionality

Adaptive Compositionality

Use more than one composition functions and adaptively select them depending on the input vectors



Adaptive Compositionality

Use more than one composition functions and adaptively select them depending on the input vectors



Adaptive Compositionality

Use more than one composition functions and adaptively select them depending on the input vectors

$$\mathbf{v} = f\left(\sum_{h=1}^{C} P(g_h | \mathbf{v}_l, \mathbf{v}_r) g_h(\mathbf{v}_l, \mathbf{v}_r)\right)$$

$$\begin{bmatrix} P(g_1 | \mathbf{v}_l, \mathbf{v}_r) \\ \vdots \\ P(g_c | \mathbf{v}_l, \mathbf{v}_r) \end{bmatrix} = softmax \begin{pmatrix} \beta S \begin{bmatrix} \mathbf{v}_l \\ \mathbf{v}_r \end{bmatrix} \end{pmatrix}$$
The Boltzmann distribution is used to adaptively select g_h .

$$\beta = 0$$

$$P(g_h | \mathbf{v}_l, \mathbf{v}_r) = \frac{1}{C}$$

$$\begin{bmatrix} P(g_1 | \mathbf{v}_l, \mathbf{v}_r) \\ \vdots \\ P(g_c | \mathbf{v}_l, \mathbf{v}_r) \end{bmatrix} = softmax \begin{pmatrix} S \begin{bmatrix} \mathbf{v}_l \\ \mathbf{v}_r \end{bmatrix} \end{pmatrix}$$

$$P(g_h | \mathbf{v}_l, \mathbf{v}_r) = \begin{cases} 1, maxmum score \\ 0, otherwise \end{cases}$$
Avg-AdaMC
Weighted-AdaMC
Max-AdaMC

Objective Function

Minimize the cross-entropy error

- Target vector $t_j = [0 ... 1 ... 0]$
- Predicted distribution $y_j = [0.07 \dots 0.69 \dots 0.15]$

$$\min_{\Theta} E(\Theta) = -\sum_{i} \sum_{j} t_{j}^{i} \log y_{j}^{i} + \sum_{\theta \in \Theta} \lambda_{\theta} \|\theta\|_{2}^{2}$$

AdaGrad (Duchi, Hazan, and Singer 2011)

$$\begin{aligned} \theta_t &= \theta_{t-1} - \eta \frac{1}{\sqrt{G_t}} \frac{\partial E}{\partial \theta} \bigg|_{\theta = \theta_{t-1}} \\ G_t &= G_{t-1} + \left(\frac{\partial E}{\partial \theta} \bigg|_{\theta = \theta_{t-1}} \right)^2 \end{aligned}$$



Parameter Estimation

- Back-propagation algorithm: $\delta_m^{i\leftarrow r}$ =

• Classification: $\frac{\partial E}{\partial \boldsymbol{U}_{mn}} = \sum_{i} [\boldsymbol{v}_{n}^{i}(\boldsymbol{y}_{m}^{i} - \boldsymbol{t}_{m}^{i})]$

$$= \begin{cases} \sum_{k} (\boldsymbol{y}_{m}^{i} - \boldsymbol{t}_{m}^{i}) \boldsymbol{U}_{mk} f'(\boldsymbol{a}_{m}^{i}), r = i \\ \sum_{k} \boldsymbol{\delta}_{m}^{par(i)\leftarrow r} \frac{\partial \boldsymbol{a}_{k}^{par(i)}}{\partial \boldsymbol{v}_{m}^{i}} f'(\boldsymbol{a}_{m}^{i}), r \in anc(i) \end{cases}$$

• Composition selection:
$$\frac{\partial E}{\partial S_{mn}} = \begin{cases} \sum_{i} \sum_{r \in bp(i)} \sum_{k} \delta_{k}^{i \leftarrow r} \sum_{h} a_{k}^{i,g_{h}} x_{n}^{i} \beta P(g_{h} | v_{l}^{i}, v_{r}^{i}) (P(g_{h} | v_{l}^{i}, v_{r}^{i}) - 1), h = m \\ \sum_{i} \sum_{r \in bp(i)} \sum_{k} \delta_{k}^{i \leftarrow r} \sum_{h} a_{k}^{i,g_{h}} x_{n}^{i} \beta P(g_{h} | v_{l}^{i}, v_{r}^{i}) P(g_{m} | v_{l}^{i}, v_{r}^{i}), h \neq m \end{cases}$$

• Linear composition:
$$\frac{\partial E}{\partial W_{mn}} = \sum_i \sum_{r \in bp(i)} \delta_m^{i \leftarrow r} \mathbf{x}_n^i P(g_h | \mathbf{v}_l^i, \mathbf{v}_r^i)$$

• Tensor Composition:
$$\frac{\partial E}{\partial v_{mn}^{h[d]}} = \sum_i \sum_{r \in bp(i)} \delta_d^{i \leftarrow r} x_m^i x_n^i P(g_h | v_l^i, v_r^i)$$

• Word Embedding:
$$\frac{\partial E}{\partial L_d^w} = \sum_{[i]=w} \sum_{r \in bp(i)} \delta_d^{i \leftarrow r}$$

Stanford Sentiment Treebank

- 10,662 critic reviews in Rotten Tomatoes
- 215,154 phrases from results of Stanford Parser
- The workers in Amazon Mechanical Turk annotate polarity levels for all these phrases
- The sentiment scales are merged to five categories (very negative, negative, neutral, positive, very positive)



Method	Fine-grained	Pos./Neg.
SVM	40.7	79.4
MNB	41.0	81.8
bi-MNB	41.9	83.1
VecAvg	32.7	80.1
MV-RNN	44.4	82.9
RNN	43.2	82.4
Avg-AdaMC-RNN	43.4	84.9
Max-AdaMC-RNN	43.8	85.6
Weighted-AdaMC-RNN	45.4	86.5
AdaMC-RNN	45.8	87.1
RNTN	45.7	85.4
Avg-AdaMC-RNTN	45.7	86.3
Max-AdaMC-RNTN	45.6	86.6
Weighted-AdaMC-RNTN	46.3	88.4
AdaMC-RNTN	<u>46.7</u>	<u>88.5</u>

Results of evaluation on the Sentiment Treebank. The top three methods are in bold. Our methods achieve best performances when \beta is set to 2.



Figure 3: The curve shows the accuracy for root nodes as $\beta = 0, 2^0, 2^1, \dots, 2^6$ increases. AdaMC-RNN and AdaMC-RNTN achieve the best results at $\beta = 2^1$.

Vector Representations

Word/Phrase	Neighboring Words/Phrases in the Vector Space
good	cool, fantasy, classic, watchable, attractive
boring	dull, bad, disappointing, horrible, annoying
ingenious	extraordinary, inspirational, imaginative, thoughtful, creative
soundtrack	execution, animation, cast, colors, scene
good actors	good ideas, good acting, good looks, good sense, great cast
thought-provoking film	beautiful film, engaging film, lovely film, remarkable film, riveting story
painfully bad	how bad, too bad, really bad, so bad, very bad
not a good movie	isn't much fun, isn't very funny, nothing new, isn't as funny of clichés



Composition Pairs in the Composition Space

- For the composition pair (v_l, v_r) , we use the distribution of the composition functions $\begin{bmatrix}
 P(g_1 | v_l, v_r) \\
 \vdots \\
 P(g_c | v_l, v_r)
 \end{bmatrix}$ to query its neighboring pairs
 - **Composition Pair Neighboring Composition Pairs** really bad very bad / only dull / much bad / extremely bad / (all that) bad (is n't) (is n't) (painfully bad) / not mean-spirited / not (too slow) / not well-acted / (necessarily bad) (have otherwise) (been bland) great (cinematic innovation) / great subject / great performance great (Broadway play) / energetic entertainment / great (comedy filmmaker) (arty and) jazzy (Smart and) fun / (verve and) fun / (unique and) entertaining / (gentle and) engrossing / (warmth and) humor

Visualization: Composition Pairs

 $\begin{bmatrix} P(g_1 | \boldsymbol{v}_l, \boldsymbol{v}_r) \\ \vdots \\ P(g_C | \boldsymbol{v}_l, \boldsymbol{v}_r) \end{bmatrix}$





- Best films
- Riveting story
- Solid cast
- Talented director
- Gorgeous visuals



- Really good
- Quite funny
- Damn fine
- Very good
- Particularly funny



- Is never dull
- Not smart
- Not a good movie
- Is n't much fun
- Wo n't be disappointed



Roberto Alagna

Pearl Harbor

Elizabeth Hurley

Diane Lane

Pauly Shore

Future Work

- Use AdaMC for the other NLP tasks
- Utilize external information to adaptively select the composition functions
 - Part-of-speech tags
 - Syntactic parsing results
- Mix different composition types together
 - Linear combination approach (RNN)
 - Tensor-based approach (RNTN)
 - Multiplication approach



Microsoft Research 微软亚洲研究院

THANKS!

