Learning Sense-specific Word Embeddings and Applying for Semantic Hierarchies

Wanxiang Che HIT-SCIR

Joint work with Jiang Guo and Ruiji Fu

2014-5-16

Word Embeddings

Definition

- Compact, Real-valued, Low-dimensional vector representations
- Example
 - *Emb(*计算机*)*=[0.12, 0.26, ..., 0.09]
 - *Emb(*电 脑)=[0.14, 0.14, ..., 0.10]

e.g. 50*dim*

- Learning Architectures
 - Context-Predicting Models
 - Neural Network Language Models (NNLM)
 - Skip-Gram Model
 - Context-Counting Models
 - LSA, CCA, etc.

Word and Sense Mapping



Learning Sense-specific Word Embedding by Exploiting Bilingual Resources

Approach

- Represent words with sense-specific embeddings
 - Word sense induction using a bilingual approach
 - Train embeddings on the sense-tagged corpus

Applications

- Word similarity evaluation of polysemous words
- Incorporating sense-specific embeddings to sequence labeling tasks (e.g. Named Entity Recognition)





Method --- Learning

- Word Alignment
- Extract the translation words via bidirectional translation probability
- Cluster the translation words using their embeddings
 - Affinity Propagation (AP) clustering is used, for automatically determine the cluster number
- Cross-lingual Sense Projection
 - Tag the words in source language with its sense cluster
- Train embeddings using RNNLM (recurrent NNLM)



Recurrent Neural Network Language Model

■ Modeling the conditional probability ■ P(w(t+1)|w(t), h(t-1))

• Compute: $\mathbf{h}(t) = f(\mathbf{U}\mathbf{w}(t) + \mathbf{W}\mathbf{h}(t-1))$ $\mathbf{y}(t) = g(\mathbf{V}\mathbf{h}(t))$

U is the Embedding Matrix



Related Work

- □ Huang et al. (2012), Reisinger and Mooney (2010)
 - Learning Multiple-prototype Word Embeddings
 - Kembeddings for each word
 - Problem
 - Ignoring the fact that the number of senses of different words are varied.

Word Similarity Evaluation

A manually constructed Polysemous Word Similarity Dataset

Measurement

- Spearman Correlation
- Kendall Correlation
- Quantitative Evaluation

System	MaxSim		AvgSim		
System	$\rho \times 100$	au imes 100	$\rho \times 100$	au imes 100	
Ours	55.4	40.9	49.3	35.2	
SingleEmb	42.8	30.6	42.8	30.6	
Multi-prototype	40.7	29.1	38.3	27.4	

Word Similarity Evaluation

Qualitative Evaluation: K-nearest neighbors

Word	Nearest Neighboors
制服 _{uniform}	穿着 _{dress} ,警服 _{policeman uniform}
制服 _{subdue}	打败 _{defeat} , 击败 _{beat} , 征服 _{conquer}
花 _{spend}	花费 _{cost} ,节省 _{save} ,剩下 _{rest}
花 _{flower}	菜 _{greens} , 叶 _{leaf} , 果实 _{fruit}
法 _{law}	法令ordinance,法案bill,法规rule
法 _{method}	概念 _{concept} ,方案 _{scheme}
法 $_{French}$	德 _{Germany} , 俄 _{Russia} , 英 _{Britain}
领导 _{lead}	监督 _{supervise} ,决策 _{decision}
领导 _{leader}	主管 _{chief} ,上司 _{boss}

Method --- Application

Semi-supervised NER

- Feed the word embeddings as additional features to a supervised learning model
- Problem: discriminate the sense of words in NER data
- Solution:
 - A novel Word Sense Disambiguation algorithm based on the RNNLM trained previously

Word Sense Disambiguation based on RNNLM

Single-step decision

 $= \operatorname{argmax} P(w(t+1) \downarrow s \uparrow * | w(t) \downarrow s \uparrow * , h(t-1))$



- Greedy decoding
- Beam-search decoding

NER Evaluation

Model: Conditional Random Fields

Result:

System	Р	R	F
Baseline	93.27	81.46	86.97
+SingleEmb	93.55	82.32	87.58
+SenseEmb (greedy)	93.38	83.56	88.20
+SenseEmb (beam search)	93.59	84.05	88.56

NER Evaluation

Analysis

Per-token accuracy of polysemous words and monosemous words, respectively



Learning Semantic Hierarchies via Word Embeddings

Semantic Hierarchies

Learning Semantic Hierarchies via Word Embeddings

- \square car \rightarrow automotive
 - hypernym: automotive
 - hyponym: car



- manually-built semantic hierarchies
 - WordNet
 - HowNet
 - CilinE (Tongyi Cilin Extended version)

Previous Work

Pattern-based method

- e.g. "such NP1 as NP2"
 - Hearst (1992); Snow et al. (2005)

Pattern	Translation		
₩ 是[一个 一种] h	w is a [a kind of] h		
w[、] 等 h	w[,] and other h		
h[,]叫[做]w	h[,] called w		
h [,] [像]如 w	h[,] such as w		
h[,] 特别是 w	h[,] especially w		

- Methods based on web mining
 - assuming that the hypernyms of an entity co-occur with it frequently
 - extracting hypernym candidates from multiple sources and learning to rank
 - Fu et al. (2013)

Word Embeddings

Learning Semantic Hierarchies via Word Embeddings

$$v(\text{king}) - v(\text{queen}) \approx v(\text{man}) - v(\text{woman})$$



Mikolov et al. (2013a)

Motivation

Does the embedding offset work well in hypernymhyponym relations?

No.	Examples
1	$v(虾) - v($ 対虾 $) \approx v($ 鱼 $) - v($ 金鱼 $)$
	$v(\text{shrimp}) - v(\text{prawn}) \approx v(\text{fish}) - v(\text{gold fish})$
2	$v($ 工人 $) - v($ 木匠 $) \approx v($ 演员 $) - v($ 小丑 $)$
	$v(\text{laborer}) - v(\text{carpenter}) \approx v(\text{actor}) - v(\text{clown})$
3	v(工人) – v(木匠) ≉ v(鱼) – v(金鱼)
	$v(\text{laborer}) - v(\text{carpenter}) \not\approx v(\text{fish}) - v(\text{gold fish})$

Motivation

Clusters of the vector offsets in training data





- Word Embedding Training
- Projection Learning
- □ is-a Relation Identification

Word Embedding Training

Skip-gram



Mikolov et al. (2013b)

Projection Learning

□ A uniform Linear Projection

Given a word x and its hypernym y, there exists a matrix Φ so that $y = \Phi x$.

$$\Phi^* = \underset{\Phi}{\operatorname{arg\,min}} \frac{1}{N} \sum_{(x,y)} \parallel \Phi x - y \parallel^2$$

Projection Learning

Piecewise Linear Projections

clustering y - x



Iearning a separate projection for each cluster

$$\Phi_k^* = \underset{\Phi_k}{\operatorname{arg\,min}} \frac{1}{N_k} \sum_{(x,y)\in C_k} \| \Phi_k x - y \|^2$$

Projection Learning

Training data



is-a Relation Identification

□ Given two words x and y

If y is determined as a hypernym of x, either of the two conditions must be satisfied.



$$d(\Phi_k x, y) = \parallel \Phi_k x - y \parallel^2 < \delta$$

Condition 2:

 $x \xrightarrow{H} z$ and $z \xrightarrow{H} y$

is-a Relation Identification

Hierarchy Condition (DAG)

•
$$\forall x, y \in L : x \xrightarrow{H} y \Rightarrow \neg(y \xrightarrow{H} x)$$

• $\forall x, y, z \in L : (x \xrightarrow{H} z \land z \xrightarrow{H} y) \Rightarrow x \xrightarrow{H} y$



Experimental Data

Word embedding training

corpus from Baidubaike

- ~30 million sentences (~780 million words)
- Projection learning
 - CilinE
 - 15,247 is-a pairs

Experimental Data

□ For evaluation



Delation	# of word pairs		
Kelation	Dev.	Test	
hypernym-hyponym	312	1,079	
hyponym-hypernym*	312	1,079	
unrelated	1,044	3,250	
Total	1,668	5,408	

Fu et al. (2013)

Results and Analysis

Comparison with existing methods

	P(%)	R (%)	F (%)	-
\mathbf{M}_{CilinE}	98.21	50.88	67.03	-
$M_{Wiki+CilinE}$	92.41	60.61	73.20	Suchanek et al. (2008)
$\mathbf{M}_{Pattern}$	97.47	21.41	35.11	Hearst (1992)
M_{Snow}	60.88	25.67	36.11	Snow et al. (2005)
$M_{balApinc}$	54.96	53.38	54.16	Kotlerman et al. (2010)
M_{invCL}	49.63	62.84	55.46	Lenci and Benotto (2012)
\mathbf{M}_{Fu}	71.64	52.92	60.87	Fu et al. (2013)
M _{offset}	59.26	63.19	61.16	-
M_{Emb}	80.54	67.99	73.74	
$M_{Emb+CilinE}$	80.59	72.42	76.29	
$M_{\it Emb+Wiki+CilinE}$	79.78	80.81	80.29	

Results and Analysis



Demo

□ <u>http://www.bigcilin.com</u>



THANKS & QA