

# LIMSI: Learning Semantic Similarity by Selecting Random Word Subsets

Artem Sokolov  
LIMSI-CNRS Orsay, France  
gposhta@gmail.com



## HIGHLIGHTS

### Semantic Textual Similarity learning between sentences

- degrees of semantic equivalence:  
5 (exactly the same meaning) .. → .. 0 (nothing in common)
- task:
  - given 2 sentences
  - predict the degree of semantic equivalence  $\in [0, 5]$
  - evaluation: Pearson's correlation
- contributions:
  - ★ improved Random Indexing (RI):
    - \* selection of informative projection directions
    - \* learning-to-rank + boosting approach
  - ★ no special preprocessing, no linguistic resources
  - ★ in top-25% among 89 participants (SemEval'12/Task6)

## SYMMETRIC VECTOR TRANSFORMATIONS

Combine 2 RI sentence vectors  $v(s_1), v(s_2)$  in a **symmetric** manner:

- ‘sumdiff’:  $\bar{x} = (v(s_1) + v(s_2), \text{sgn}(v_1(s_1) - v_1(s_2))(v(s_1) - v(s_2)))$
- ‘concat’:  $\bar{x} = (v(s_1), v(s_2))$ , and  $\bar{x}' = (v(s_2), v(s_1))$
- ‘product’:  $x_i = v_i(s_1) \cdot v_i(s_2)$
- ‘crossprod’:  $x_{ij} = v_i(s_1) \cdot v_j(s_2)$
- ‘crossdiff’:  $x_{ij} = v_i(s_1) - v_j(s_2)$
- ‘absdiff’:  $x_i = |v_i(s_1) - v_i(s_2)|$ .

## PERFORMANCE ON TRAINING DATA

	learner	transform	correlation	$1\sigma$
baseline	pure RI, cos	-	0.264	0.005
	logistic reg.	-	0.508	0.041
	logistic reg.	concat	0.537	0.052
boosting	RankBoost	sumdiff	0.685	0.027
		product	0.663	0.018
		crossprod	0.648	0.028
		crossdiff	0.643	0.023
		concat	0.625	0.025
		absdiff	0.602	0.021
RtRank	sumdiff	0.730	0.020	
	product	0.721	0.023	

- selected transformations for the submission **in bold**.
- optimal parameters for RankBoost & RtRank found with cross-validation on 5 folds

## CLASSIC RANDOM INDEXING [LANDAUER AND DUMAIS1997]

Distributional hypothesis: words with similar meanings tend to appear in similar contexts

### map words $w$ to context vectors $v(w)$ :

- 1 initialize with ternary sparse vectors
- 2 sum over contexts
- 3-4 normalize for frequency and sentence length
  - ✓ similar contexts ⇒ similar vectors

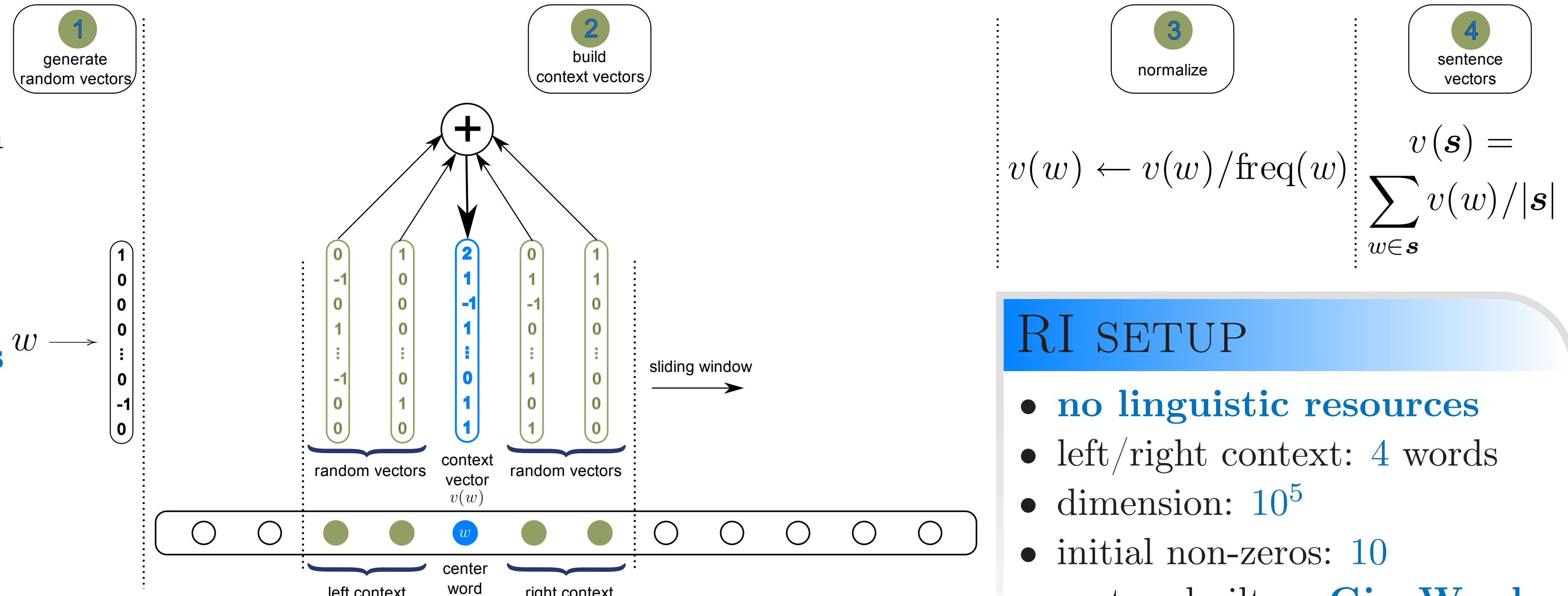
### ≈ dim. reduction by random projections

- ★  $\ell_2$ -norm between bags-of-words of contexts
- ★ ternarity turns on/off random word sets

### drawbacks:

- ✗ agnostic of the modeled semantic relation
- ✗ trial-and-error tuning to specific task
- ✗ all (even useless) random sets have effect

### adaption to different tasks needed!



## RI SETUP

- no linguistic resources
- left/right context: 4 words
- dimension:  $10^5$
- initial non-zeros: 10
- vectors built on **GigaWord** (no punctuation, lowercased)

## CONTRIB.: SELECTING WORD SUBSETS WITH BOOSTING TO RANK BY SEMANTIC SIMILARITY

### Idea:

- ★ exploit ordering of sentence pairs under  $y^i$
- ★ use boosting to select important elements of  $\bar{x}$  ⇔ select random word sets relevant to the similarity modeled

### Given:

- training sentence pair:  $s_1^i, s_2^i$  & similarity:  $y^i$

### Want to Learn:

- similarity score  $H(\bar{x}) = \sum_{i=1}^T \alpha_t h_t(\bar{x})$ , where
  - $\bar{x}$  – some symmetric transform of sentence vectors  $v(s_1), v(s_2)$  (see frame to the left)
  - $h_t(\bar{x})$  – **simple (weak) functions** of  $\bar{x}$ ,  $\alpha_t$  – weights
  - the more similar  $v(s_1)$  is to  $v(s_2)$ , the higher should be  $H(\bar{x})$

## CONCLUSION & EXPERIMENTAL RESULTS

- ★ **RI extention**: selects word sets relevant for semantics learning
- ★ learns semantic order, not absolute values
- ★ **promising results**, despite ignoring preprocessing & linguistics
- underperforms on SMT sets: possible incorrect lexical choice in MT output + should be trained separately

learner	transform	train $\pm \sigma$	test	final rank	MSRpar	MSRvid	SMTeur	OnWN	SMTnews
RankBoost	product	0.748 $\pm$ 0.017	0.6392	32	0.3948	0.6597	0.0143	0.4157	0.2889
	sumdiff	0.735 $\pm$ 0.016	0.6196	45	0.4295	0.5724	0.2842	0.3989	0.2575
RtRank	product	0.784 $\pm$ 0.017	<b>0.6789</b>	<b>22</b>	0.4848	0.6636	0.0934	0.3706	0.2455
	sumdiff	0.763 $\pm$ 0.014							

### Ranking losses & Weak functions $h_t$ :

- **RankBoost** [Freund et al.2003]
  - loss:  $\sum_{(s_1^i, s_2^i), (s_1^j, s_2^j): y^i < y^j} P(i, j) \llbracket H(\bar{x}^i) \geq H(\bar{x}^j) \rrbracket$
  - weak functions:  $h(x; \theta, k) = \llbracket x_k > \theta \rrbracket$  (decision stump)
  - weights  $P(i, j) = y^j - y^i$
- **RtRank** [Zheng et al.2007, Mohan et al.2011]
  - loss:  $\sum_{(s_1^i, s_2^i), (s_1^j, s_2^j): y^i < y^j} (\max\{0, H(\bar{x}^i) - H(\bar{x}^j)\})^2$
  - weak functions: decision trees of depth 4
  - Well-known algorithms, implementations available online