# Learning to Reweight Terms with Distributed Representations

Guoqing Zheng, Jamie Callan

School of Computer Science Carnegie Mellon University

August 12, 2015

### Outline

Goal:

Assign weights to query terms for better retrieval results

- Motivation
- Method
- Evaluation
- Conclusions

### Motivation

Assigning weights to terms is not a new problem

$$score(q, d) = \sum_{t \in q} s(t, q, C) f(t, d)$$
 (1)

Brief summary of existing representative methods:

Туре	Method	Performance	Features	Model
Unsupervised	Constant	bad	-	-
Olisupervised	IDF	average	-	-
Supervised	LtR-methods	good	manual,	non-convex
(Rank measure)	(WSD)	good	external	non-convex
Supervised (Term recall)	Zhao et al	average	manual	requires SVD

 Can we have something that is directly term weight supervised, can lead to good retrieval performance without manual / expensive feature construction?



# Probably yes?

- What might be a good choice of term weights? Term recall
- The probability of term t of query q appears in relevant documents  $R_q$ , can be estimated as  $\frac{|R_{q,t}|}{|R_q|}$  if relevance judgments are available
- True term recall proves to be highly effective as shown by Zhao et al[3, 2]
- But which features are good to predict that?

### Distributed word vectors

 Distributed word vectors learnt from a large corpora can be used to effectively measure semantic similarity

$$vec(\mathsf{King}) - vec(\mathsf{Man}) \approx vec(\mathsf{Queen}) - vec(\mathsf{Woman})$$

 Another important characteristic of word vector is the addition property, i.e., it's often semantically meaningful to construct word vectors for phrases from those of individual terms, such as representing

$$vec(\mathsf{Chinese} \; \mathsf{river}) \approx vec(\mathsf{Chinese}) + vec(\mathsf{river})$$



### Idea in this work

- Intuitively, term recall measures the relative importance of a term w.r.t the whole query in determining document-query relevance
- With the addition property of word vectors, we can construct feature vectors for terms from their word vector representations, which resembles the above intuition of term recall
- For a term t in query q, we construct the feature vector x(t) given the word vector representation w(t) as

$$x(t) \equiv vec(t) - vec(q)$$

- where vec(q) is the average word vector of all the terms in q.
- The above defined feature vector is essentially the offset vector of a term to the center of the entire query



### Feature vector illustration

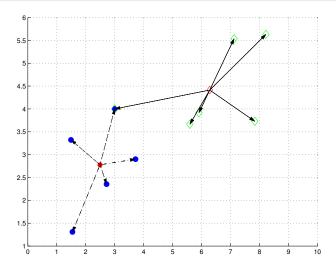


Figure: Demonstration of transforming global 2-dimensional word vectors into feature vectors local to the query.

# Term recall learning with feature vectors

- Model term recall as a function of the feature vectors defined
- Preprocessing: transform term recall r from [0,1] to  $y \in \mathbb{R}$  by a logit function y = logit(r)
- ullet Models the transformed term recall y as a linear function of the feature vectors x
- Employ  $\ell_1$  norm regularization in learning ( $\ell_2$  norm also tried)
- LASSO optimization:

$$\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{n_i} (y_{ij} - \beta^{\top} x_{ij})^2 + \lambda \|\beta\|_1$$
 (2)

 $oldsymbol{\circ}$  Predict term recall for testing query term with feature vector  $oldsymbol{x}_*$  as

$$\widehat{\mathbb{P}(t_*|R)} = \operatorname{sigmoid}\left(\widehat{\boldsymbol{\beta}}^{\top} \boldsymbol{x}_*\right) = \frac{\exp\{\widehat{\boldsymbol{\beta}}^{\top} \boldsymbol{x}_*\}}{1 + \exp\{\widehat{\boldsymbol{\beta}}^{\top} \boldsymbol{x}_*\}}$$
(3)

# Incorporating predicted weights to query models (I)

#### Base query models:

```
Bag-of-Words (BOW):
 apple pie recipe

    Sequential Dependency model (SD):

 #weight(
 0.8 #combine(apple pie recipe)
 0.1 #combine( #1(apple pie)
                 #1(pie recipe) )
 0.1 #combine(
                 #uw8(apple pie)
                 #uw8(pie recipe) ) )
```

# Incorporating predicted weights to query models (II)

### Term recall enhanced query models:

```
DeepTR-BOW:
  #weight( \widehat{\mathbb{P}}(\mathsf{apple}|R) apple
                 \mathbb{P}(\mathsf{pie}|R) pie
                  \widehat{\mathbb{P}}(\text{recipe}|R) recipe )
DeepTR-SD:
   #weight(
  0.8 #weight( \widehat{\mathbb{P}}(\mathsf{apple}|R) apple
                        \widehat{\mathbb{P}}(\mathsf{pie}|R) pie
                        \widehat{\mathbb{P}}(\mathtt{recipe}|R) recipe )
   0.1 #combine(#1(apple pie)
                        #1(pie recipe) )
   0.1 #combine(
                         #uw8(apple pie)
                         #uw8(pie recipe) ) )
```

### **Evaluation**

Data sets

Collection	# Docs	# Words	TREC Topics
ROBUST04	528K	253M	351-450, 601-700
WT10g	1,692K	1,076M	451-550
GOV2	25,205K	24,007M	701-850
ClueWeb09B	29,038K	23,890M	1-200

- Baseline query models: Bag of Words (BOW), Sequential Dependency model (SD), Weighted Sequential Dependency model (WSD) [1]
- Retrieval models: Language model and BM25
- Evaluation metrics: P@10, MAP for all collections, NDCG@10 and ERR@20 for collections with graded relevance judgments

### **Evaluation**

#### Word vector source:

- trained from all above collections with dimensions ranging from 100, 300, 500
- pre-trained by Google trained from Google News (100 billion words) with dimension 300

# Experimental results

# Retrieval results of DeepTR-BOW with language model (dimension=300)

(				
Query Model	ROBUST04	WT10g	GOV2	ClueWeb09B
Query Moder	MAP	MAP	MAP	MAP
BOW	0.2512	0.1943	0.2488	0.0702
SD	0.2643	0.2032	0.2688	0.0745
WSD	0.2749	0.2260	0.2946	-
DeepTR-BOW	$0.2591^{b}$	0.2103	$0.2646^{b}$	0.0718
(Corpus)	(+3.2/-1.9)	(+8.2/+3.5)	(+6.3/-1.6)	(+2.2/-3.6)
DeepTR-BOW	$0.2650^{b}$	$0.2111^{b}$	$0.2646^{b}$	0.0741
(GOV2)	(+5.5/+0.3)	(+8.7/+3.9)	(+6.3/-1.6)	(+5.6/-0.5)
DeepTR-BOW	$0.2657^{b}$	0.2129	$0.2685^{b}$	0.0718
(ClueWeb09B)	(+5.8/+0.5)	(+9.6/+4.8)	(+7.9/-0.1)	(+2.2/-3.6)
DeepTR-BOW	$0.2673^{b}$	$0.2122^{b}$	$0.2630^{b}$	0.0732
(Google)	(+6.4/+1.2)	(+9.3/+4.5)	(+5.7/-2.2)	(+4.2/-1.8)

b: Statistically significant difference with BOW

 $s: \mathsf{Statistically}$  significant difference with  $\mathsf{SD}$ 

 $\Rightarrow$  Weighted BoW queries significantly outperforms unweigted BoW queries and are even comparable to SD queries.



# Retrieval results of DeepTR-SD with language model (dimension=300)

(4					
Query Model	ROBUST04	WT10g	GOV2	ClueWeb09B	
	MAP	MAP	MAP	MAP	
BOW	0.2512	0.1943	0.2488	0.0702	
SD	0.2643	0.2032	0.2688	0.0745	
WSD	0.2749	0.2260	0.2946	-	
DeepTR-SD	$0.2754_s^b$	$0.2182_s^b$	$0.2831_s^b$	0.0748	
(Corpus)	(+9.6/+4.2)	(+12.3/+7.4)	(+13.8/+5.3)	(+6.5/+0.5)	
DeepTR-SD	$0.2781_s^b$	$0.2223_s^b$	$0.2831_s^b$	$0.0806_s^b$	
(GOV2)	(+10.7/+5.2)	(+14.4/+9.4)	(+13.8/+5.3)	(+14.8/+8.2)	
DeepTR-SD	$0.2810_s^b$	$0.2279_s^b$	$0.2890_s^b$	0.0748	
(ClueWeb09B)	(+11.9/+6.3)	(+17.3/+12.1)	(+16.2/+7.5)	(+6.5/+0.5)	
DeepTR-SD	$0.2842_s^b$	$0.2256_s^b$	$0.2814_s^b$	$0.0780_s^b$	
(Google)	(+13.1/+7.5)	(+16.1/+11.0)	(+13.1/+4.7)	(+11.0/+4.7)	

b: Statistically significant difference with BOW

 $\boldsymbol{s}$  : Statistically significant difference with SD

 $\Rightarrow$  DeepTR-SD significantly outperforms both BoW and SD queries.



### Retrieval results of DeepTR-BOW with BM25 (dimension=300)

Query Model	ROBUST04	WT10g	GOV2	ClueWeb09B
Query Model	MAP	MAP	MAP	MAP
BOW	0.2460	0.1783	0.2334	0.0566
SD	0.2546	0.1817	0.2409	0.0600
DeepTR-BOW	$0.2566^{b}$	$0.1903^{b}$	$0.2430^{b}$	$0.0593^{b}$
(Corpus)	(+4.3/+0.8)	(+6.7/+4.7)	(+4.1/+0.9)	(+4.7/-1.2)
DeepTR-BOW	$0.2561^{b}$	$0.1910^{b}$	$0.2430^{b}$	$0.0596^{b}$
(GOV2)	(+4.1/+0.6)	(+7.1/+5.1)	(+4.1/+0.9)	(+5.3/-0.7)
DeepTR-BOW	$0.2590^{b}$	$0.1932^{b}$	$0.2462^{b}$	$0.0593^{b}$
(ClueWeb09B)	(+5.3/+1.8)	(+8.3/+6.3)	(+5.5/+2.2)	(+4.7/-1.2)
DeepTR-BOW	$0.2600^{b}$	$0.1922^{b}$	$0.2449^{b}$	0.0584
(Google)	(+5.7/+2.1)	(+7.8/+5.7)	(+4.9/+1.7)	(+3.1/-2.7)

b: Statistically significant difference with BOW s: Statistically significant difference with SD

⇒ Similar results observed for BM25 retrieval.

### Retrieval results of DeepTR-SD with BM25 (dimension=300)

Query Model	ROBUST04	WT10g	GOV2	ClueWeb09B
Query Moder	MAP	MAP	MAP	MAP
BOW	0.2460	0.1783	0.2334	0.0566
SD	0.2546	0.1817	0.2409	0.0600
DeepTR-SD	$0.2595_s^b$	$0.1944_s^b$	$0.2478_s^b$	$0.0613_s^b$
(Corpus)	(+5.5/+1.9)	(+9.0/+7.0)	(+6.1/+2.9)	(+8.2/+2.1)
DeepTR-SD	$0.2609_s^b$	$0.1941_{s}$	$0.2478_{s}^{b}$	$0.0616_s^b$
(GOV2)	(+6.1/+2.5)	(+8.8/+6.8)	(+6.1/+2.9)	(+8.8/+2.6)
DeepTR-SD	$0.2630_s^b$	$0.1951_s^b$	$0.2502_s^b$	$0.0613_s^b$
(ClueWeb09B)	(+6.9/+3.3)	(+9.4/+7.3)	(+7.2/+3.9)	(+8.2/+2.1)
DeepTR-SD	$0.2627_{s}^{b}$	$0.1938_{s}$	$0.2461^{b}$	0.0604
(Google)	(+6.8/+3.2)	(+8.7/+6.7)	(+5.4/+2.2)	(+6.8/+0.7)

b: Statistically significant difference with BOW s: Statistically significant difference with SD

⇒ Similar results observed for BM25 retrieval.

# Word vector dimensionality

vectors trained from ClueWeb09B, Language Model retrieval

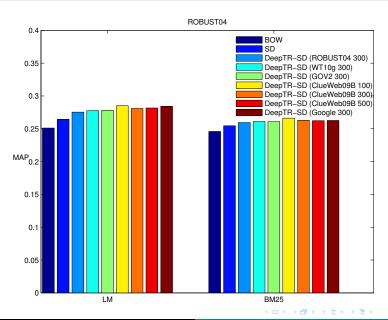
Query Model	ROBUST04	WT10g	GOV2	ClueWeb09B
Query Woder	MAP	MAP	MAP	MAP
BOW	0.2512	0.1943	0.2488	0.0702
SD	0.2643	0.2032	0.2688	0.0745
WSD	0.2749	0.2260	0.2946	-
DeepTR-BOW	$0.2681^{b}$	$0.2239^{b}$	$0.2667^{b}$	0.0727
(100)	(+6.7/+1.4)	(+15.3/+10.2)	(+7.2/-0.8)	(+3.6/-2.4)
DeepTR-BOW	$0.2657^{b}$	0.2129	$0.2685^{b}$	0.0718
(300)	(+5.8/+0.5)	(+9.6/+4.8)	(+7.9/-0.1)	(+2.2/-3.6)
DeepTR-BOW	$0.2679^{b}$	$0.2179^{b}$	$0.2655^{b}$	0.0726
(500)	(+6.6/+1.4)	(+12.1/+7.2)	(+6.7/-1.2)	(+3.4/-2.5)
DeepTR-SD	$0.2851_s^b$	$0.2373_s^b$	$0.2848_s^b$	$0.0778_s^b$
(100)	(+13.5/+7.9)	(+22.1/+16.8)	(+14.4/+5.9)	(+10.8/+4.5)
DeepTR-SD	$0.2810_s^b$	$0.2279_s^b$	$0.2890_s^b$	0.0748
(300)	(+11.9/+6.3)	(+17.3/+12.1)	(+16.2/+7.5)	(+6.5/+0.5)
DeepTR-SD	$0.2817_{s}^{b}$	$0.2306_s^b$	$0.2860_s^b$	$0.0781_{s}^{b}$
(500)	(+12.2/+6.6)	(+18.7/+13.5)	(+14.9/+6.4)	(+11.3/+4.9)

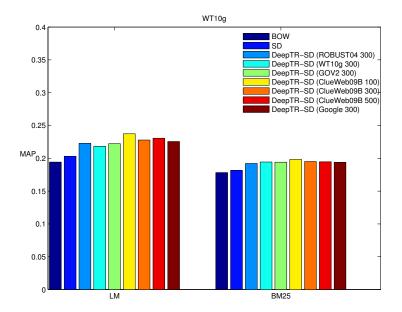
b : Statistically significant difference with BOW

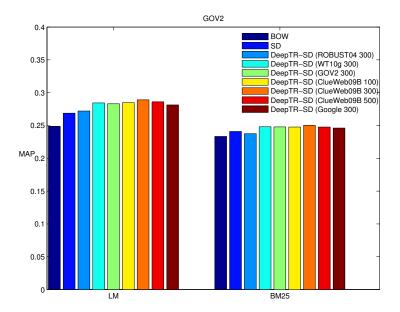
 $\boldsymbol{s}$  : Statistically significant difference with SD

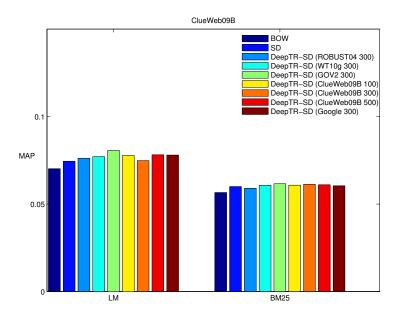
 $\Rightarrow d = 100$  works slightly better



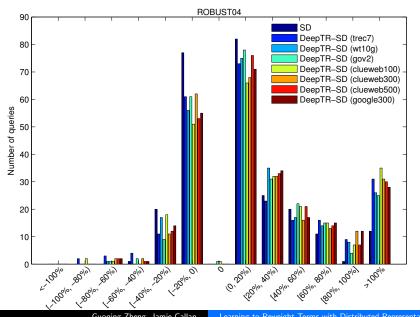


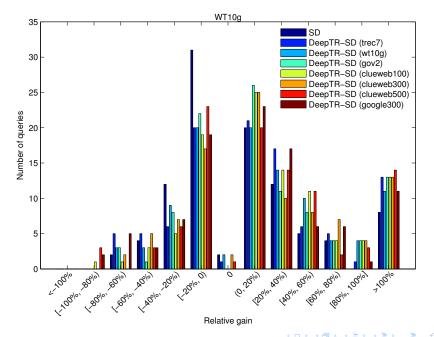


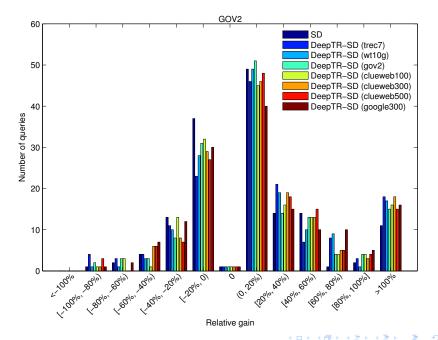


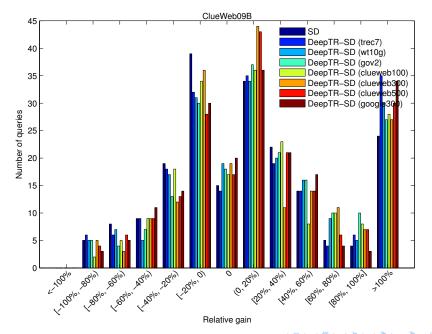


### Robustness

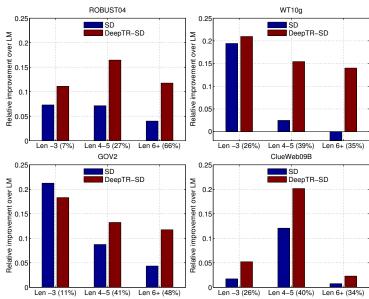








# Query length



## Graded relevance

Language Model Retrieval				
Query Model	GOV2		ClueWeb09B	
	NDCG@20	ERR@20	NDCG@20	ERR@20
BOW	0.3732	0.1493	0.1597	0.1253
SD	0.4059	0.1607	0.1694	0.1243
DeepTR-SD (GOV2)	$0.4121^{b}$	$0.1626^{b}$	$0.1743^{b}$	$0.1312_{s}$
DeepTR-SD (ClueWeb09B)	0.4146 <sup>b</sup>	$0.1614^{b}$	0.1663	0.1255
DeepTR-SD (Google)	$0.4112^{b}$	0.1600	0.1698	0.1302

	BM25 Retrieval			
Query Model	GOV2		ClueWeb09B	
	NDCG@20	ERR@20	NDCG@20	ERR@20
BOW	0.3789	0.1487	0.1188	0.1075
SD	0.3940	0.1554	0.1294	0.1059
DeepTR-SD (GOV2)	$0.4008_s^b$	$0.1590_{s}^{b}$	0.1307	0.1068
DeepTR-SD (ClueWeb09B)	$0.4033_s^b$	$0.1587_s^b$	0.1305	0.1067
DeepTR-SD (Google)	$0.4010_s^b$	$0.1586_{s}$	0.1298	0.1050

b : Statistically significant difference with BOW

 $<sup>\</sup>boldsymbol{s}\,:\,\mathsf{Statistically}$  significant difference with SD

### Conclusions and future work

- Novel (simple) framework for learning term weightes from features directly constructed from distributed word vectors
- Extensive experiments with two retrieval models on four test collections demonstrates the effectiveness

Completing the table:

Туре	Method	Performance	Features	Model
Unsupervised	Constant	bad	-	-
Olisupervised	IDF	average	-	-
Supervised	LtR-methods	good	manual,	non-convex
(Rank measure)	(WSD)	good	external	non-convex
Supervised (Term recall)	Zhao et al	average	manual	requires SVD
Supervised (Term recall)	Proposed work	good	automatic	convex, cheap

 Future directions include predicting term recalls for bigrams and proximity terms, other possible ways to construct feature vectors from distributed vectors, using word vectors for query expansion, theoretical justifications, etc.

# Acknowledgements and QA

- Reproducibility: codes, preprocessed queries, experimental setups are available
- Contact: gzheng@cs.cmu.edu
- Thanks to SIGIR for the student travel grant
- Questions?