S UNIVERSITY IN MARYLA

Mixed Membership Word Embeddings for Computational Social Science

James Foulds

University of Maryland, Baltimore County

Overview

- Word embeddings find similarities between words, leading to **improved performance** for many NLP tasks: - translation, part-of-speech tagging, chunking, NER, ...
- Allow NLP "from scratch," without feature engineering
- Typically trained in **big data** setting
- Have not yet been widely adopted for **computational social science** research due to the following limitations:
 - Target corpus of interest is often **not big data**
 - It is important for the model to be **interpretable**
- I propose a method for training **interpretable word** embeddings without big data, for computational social science, leveraging insights from topic models

Connections to Topic Models, and Mixed Membership Extension to the Skip-Gram

Skip-gram corresponds to a supervised naïve Bayes topic model, up to its parameterization via embeddings

I propose **topic model** and **mixed membership** variants

- Mixed membership models provide parameter sharing
- Can use **fewer vectors than words** for small data, while retaining substantial representational power

	Skip-gram	Skip-gram topic model	Skip-gram topic model		
	For each word in the corpus w_i	arameterize For each word in the corpus w_i			
Naive Bayes (Mikolov et a	For each word $w_c \in context(i)$	For each word $w_c \in context(i)$	mixed membership extension		
	al., 2013) Draw w_c via	Draw w_c via		1	
	$p(w_c w_i) \propto exp(v'_{w_c} {}^{T} v_{w_i} + b_{w_c})$) $p(w_c w_i) = \text{Discrete}(\phi^{(w_i)})$			
	For each word in the corpus w_i	For each word in the corpus w_i			
Mixed membership	Draw a topic $z_i \sim \text{Discrete}(\theta^{(w_i)})$	Draw a topic $z_i \sim \text{Discrete}(\theta^{(w_i)})$			
	For each word $w_c \in context(i)$	For each word $w_c \in context(i)$			

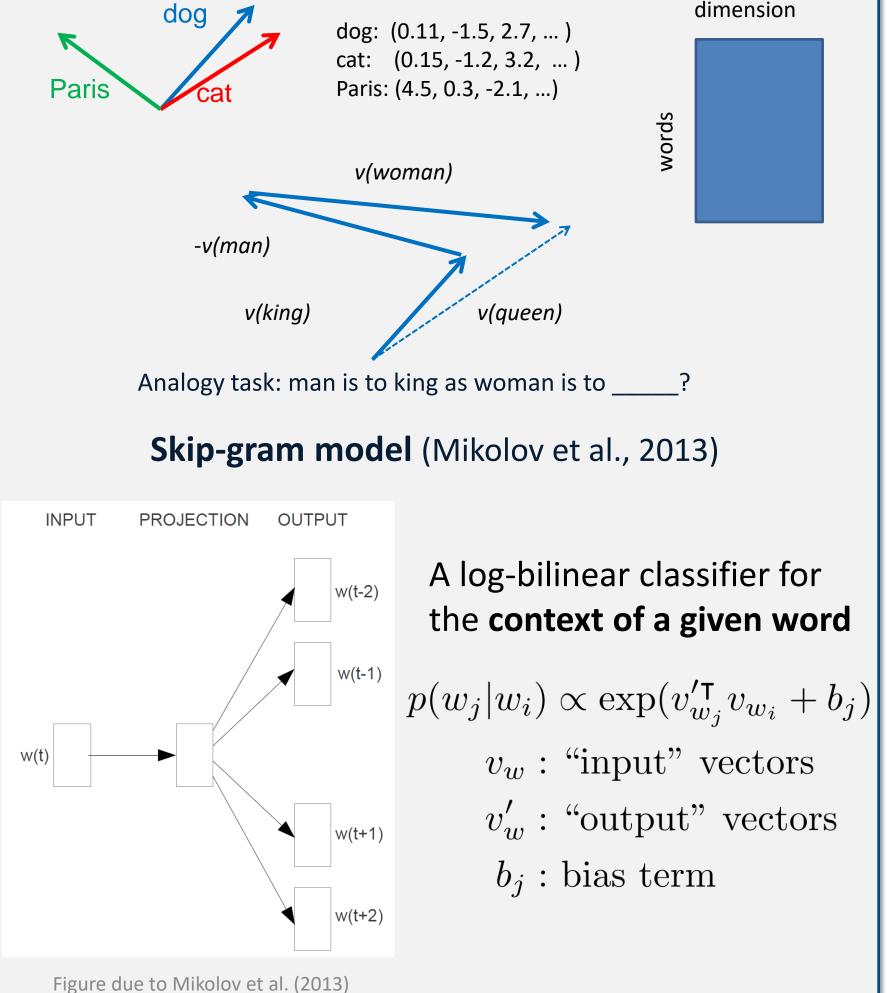


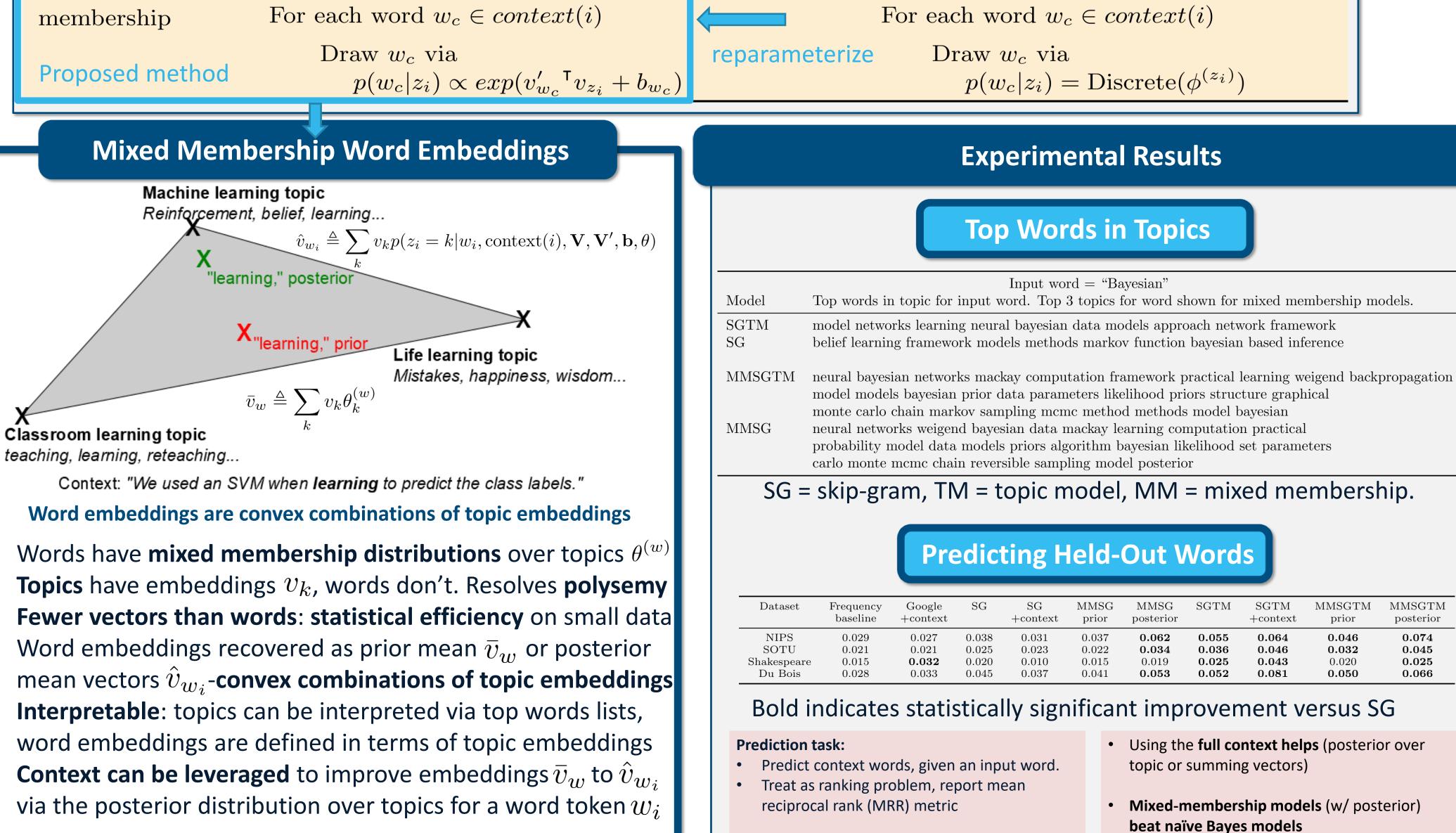
Contributions

- Interpretable, statistically efficient embedding model
- Key insight: Mixed membership representation for parameter sharing while retaining model flexibility
- **Efficient training algorithm**, using recent advances from both topic models and word embeddings:
- Metropolis-Hastings-Walker (Li et al., 2014)
- Noise-contrastive estimation (Gutmann and Hyvarinen, 2010, 2012) Proposed training algorithm is (amortized) sublinear time in the vocabulary size and number of topics
- Extensive quantitative experimental results
- Computational social science case studies
- **Practical recommendations** and insights based on these results, especially on the use of generic big data embeddings, which is a very common practice in NLP

Background: Word Embeddings

Represent **dictionary words** with **vectors**. Similar words have similar vectors.





- Simple model, scales easily to large data sets
 - Can beat deep neural network models

Inference for MM Skip-Gram Topic Model

Bayesian inference via collapsed Gibbs sampling

 $p(z_i = k|\cdot) \propto \left(n_k^{(w_i)\neg i} + \alpha_k\right) \prod_{c=1}^{|\text{context}(i)|} \frac{n_{w_c^{(i)}}^{(k)\neg i} + \beta_{w_c^{(i)} + n_{w_c^{(i)}}^{(i,c)}}}{n^{(k)\neg i} + \sum_{w'} \beta_{w'} + c - 1}$

- Scale to many topics: Metropolis-Hastings-Walker
- Alias table data structure, amortized O(1) sampling
- "Mixture of experts" proposal, alias tables for words

 $c \sim \text{Uniform}(|\text{context}(w_i)|), q_{w_c}(k) \propto \frac{n_{w_c}^{(k)} + \beta_{w_c}}{n^{(k)} + \sum_{w'} \beta_{w'}}$

Simulated annealing to escape early local maxima

Amortized Sublinear Time Training for MM Skip-Gram

- Online EM impractical O(KV) updates
- Key insight: MMSG topic model equivalent to word embedding model (*up to the Dirichlet prior*)
 - Pre-solve E-step via topic model CGS MHW algorithm
 - Apply noise-contrastive estimation to solve M-step

data vectors (except for Shakespeare dataset) • Topic models beat embedding models

Training on target corpus beats generic big-

Downstream Tasks: Classification and Regression

Dataset	#Classes	#Topics	Tf-idf	Google	MMSG	SG	MMSGTM	SG+MMSG	SG+MMSG+Google
20 Newsgroups Reuters-150 Ohsumed	$20 \\ 150 \\ 23$	$200 \\ 500 \\ 500$	$83.33 \\ 73.04 \\ 43.07$	$52.50 \\ 53.65 \\ 20.56$	$55.58 \\ 65.26 \\ 31.82$	$59.50 \\ 69.53 \\ 37.57$	$64.08 \\ 66.97 \\ 32.41$	$\begin{array}{c} 66.55 \\ 70.63 \\ 39.53 \end{array}$	$72.53 \\ 71.20 \\ 40.27$
SOTU (RMSE)	Regression	500	19.57	8.64	12.73	10.57	21.88	9.94	8.15

- Document categorization (classification accuracy, larger is better), and predicting the year of SOTU addresses (RMSE, smaller is better).
- **Target corpus beats generic big-data vectors** (except for SOTU, which is very small)
- Skip-gram beats MMSG for classification/regression loss of granularity
- But, concatenating the different vectors improves performance over individual embeddings
 - MMSG, SG, generic Google vectors learn complementary information

Vector Composition in Topic Space

Nearest topic after composition of mean vectors for words

object + recognitioncharacter + recognitionspeech + recognitioncomputer + visioncomputer + sciencebias + variance covariance + variance

objects visual object recognition model recognition segmentation character speech recognition hmm system hybrid computer vision ieee image pattern university science colorado department error training set data performance gaussian distribution model matrix

Data Visualization: Document, Topic, and Author Embeddings on State of the Union Addresses and NIPS Articles (t-SNE Projections)

State of the Union Addresses

Democrats (blue), liberal topics

NIPS Documents

neural, networks, ieee, systems, hopfield

