# Mixed Membership Word Embeddings
# for Computational Social Science

**James R. Foulds**

Department of Information Systems, University of Maryland, Baltimore County.

## Abstract

Word embeddings improve the performance of NLP systems by revealing the hidden structural relationships between words. Despite their success in many applications, word embeddings have seen very little use in computational social science NLP tasks, presumably due to their reliance on big data, and to a lack of interpretability. I propose a probabilistic model-based word embedding method which can recover *interpretable embeddings, without big data.* The key insight is to leverage *mixed membership* modeling, in which global representations are shared, but individual entities (i.e. dictionary words) are free to use these representations to uniquely differing degrees. I show how to train the model using a combination of state-of-the-art training techniques for word embeddings and topic models. The experimental results show an improvement in predictive language modeling of up to 63% in MRR over the skip-gram, and demonstrate that the representations are beneficial for supervised learning. I illustrate the interpretability of the models with computational social science case studies on State of the Union addresses and NIPS articles.

## 1 Introduction

Word embedding models, which learn to encode dictionary words with vector space representations, have been shown to be valuable for a variety of natural language processing (NLP) tasks such as statistical machine translation (Vaswani et al., 2013), part-of-speech tagging, chunking, and named entity recogition (Col-

lobert et al., 2011), as they provide a more nuanced representation of words than a simple indicator vector into a dictionary. These models follow a long line of research in data-driven semantic representations of text, including latent semantic analysis (Deerwester et al., 1990) and its probabilistic extensions (Hofmann, 1999a; Griffiths et al., 2007). In particular, topic models (Blei et al., 2003) have found broad applications in computational social science (Wallach, 2016; Roberts et al., 2014) and the digital humanities (Mimno, 2012), where interpretable representations reveal meaningful insights. Despite widespread success at NLP tasks, word embeddings have not yet supplanted topic models as the method of choice in computational social science applications. I speculate that this is due to two primary factors: 1) a perceived reliance on big data, and 2) a lack of interpretability. In this work, I develop new models to address both of these limitations.

Word embeddings have risen in popularity for NLP applications due to the success of models designed specifically for the big data setting. In particular, Mikolov et al. (2013a,b) showed that very simple word embedding models with high-dimensional representations can scale up to massive datasets, allowing them to outperform more sophisticated neural network language models which can process fewer documents. In this work, I offer a somewhat contrarian perspective to the currently prevailing trend of big data optimism, as exemplified by the work of Mikolov et al. (2013a,b); Collobert et al. (2011), and others, who argue that massive datasets are sufficient to allow language models to automatically resolve many challenging NLP tasks. Note that "big" datasets are not always available, particularly in computational social science NLP applications, where the data of interest are often not obtained from large scale sources such as the internet and social media, but from sources such as press releases (Grimmer, 2010), academic journals (Mimno, 2012), books (Zhu et al., 2015), and transcripts of recorded speech (Brent, 1999; Nguyen et al., 2014; Guo et al., 2015).

A standard practice in the literature is to train word embedding models on a generic large corpus such as

| NIPS | reinforcement belief learning policy algorithms Singh robot machine MDP planning |
|---|---|
| **Google News** | teaching learn learning reteaching learner_centered emergent_literacy kinesthetic_learning |

Table 1: Most similar words to "*learning*," based on word embeddings trained on NIPS articles, and on the large generic Google News corpus (Mikolov et al., 2013a,b).

Wikipedia, and use the embeddings for NLP tasks on the target dataset, cf. (Collobert et al., 2011; Mikolov et al., 2013a; Pennington et al., 2014; Kiros et al., 2015). However, as we shall see here, this standard practice might not always be effective, as the size of a dataset does not correspond to its degree of relevance for a particular analysis. Even very large corpora have idiosyncrasies that can make their embeddings invalid for other domains. For instance, suppose we would like to use word embeddings to analyze scientific articles on machine learning. In Table 1, I report the most similar words to the word "*learning*" based on word embedding models trained on two corpora. For embeddings trained on articles from the NIPS conference, the most similar words are related to *machine learning*, as desired, while for embeddings trained on the massive, generic Google News corpus, the most similar words relate to *learning and teaching in the classroom.* Evidently, domain-specific data can be important.

Even more concerningly, Bolukbasi et al. (2016) show that word embeddings can encode implicit sexist assumptions. This suggests that when trained on large generic corpora they could also encode the hegemonic worldview, which is inappropriate for studying, e.g., black female hip-hop artists' lyrics, or poetry by Syrian refugees, and could potentially lead to systematic bias against minorities, women, and people of color in NLP applications with real-world consequences, such as automatic essay grading and college admissions. In order to proactively combat these kinds of biases in large generic datasets, and to address computational social science tasks, there is a need for effective word embeddings for small datasets, so that the most relevant datasets can be used for training, even when they are small. To make word embeddings a viable alternative to topic models for applications in the social sciences, we further desire that the embeddings are semantically meaningful to human analysts.

In this paper, I introduce an interpretable word embedding model, and an associated topic model, which are designed to work well when trained on a small to medium-sized corpus of interest. The primary insight is to use a data-efficient parameter sharing scheme via *mixed membership* modeling, with inspiration from topic models. Mixed membership models provide a flexible yet efficient latent representation, in which entities are associated with shared, global representations,

but to uniquely varying degrees. I identify the skip-gram word2vec model of Mikolov et al. (2013a,b) as corresponding to a certain naive Bayes topic model, which leads to mixed membership extensions, allowing the use of *fewer vectors than words*. I show that this leads to better modeling performance without big data, as measured by predictive performance (when the context is leveraged for prediction), as well as to interpretable latent representations that are highly valuable for computational social science applications. The interpretability of the representations arises from defining embeddings for words (and hence, documents) in terms of embeddings for topics. My experiments also shed light on the relative merits of training embeddings on generic big data corpora versus domain-specific data.

## 2 Background

In this section, I provide the necessary background on word embeddings, as well as on topic models and mixed membership models. Traditional language models aim to predict words given the contexts that they are found in, thereby forming a joint probabilistic model for sequences of words in a language. Bengio et al. (2003) developed improved language models by using *distributed representations* (Hinton et al., 1986), in which words are represented by neural network synapse weights, or equivalently, vector space embeddings.

Later authors have noted that these *word embeddings* are useful for semantic representations of words, independently of whether a full joint probabilistic language model is learned, and that alternative training schemes can be beneficial for learning the embeddings. In particular, Mikolov et al. (2013a,b) proposed the *skip-gram* model, which inverts the language model prediction task and aims to *predict the context* given an input word. The skip-gram model is a log-bilinear discriminative probabilistic classifier parameterized by "input" word embedding vectors $v_{w_i}$ for the input words $w_i$, and "output" word embedding vectors $v'_{w_c}$ for context words $w_c \in \text{context}(i)$, as shown in Table 2, top-left.

Topic models such as *latent Dirichlet allocation* (LDA) (Blei et al., 2003) are another class of probabilistic language models that have been used for semantic representation (Griffiths et al., 2007). A straightforward way to model text corpora is via unsupervised multinomial naive Bayes, in which a latent cluster assignment

**James R. Foulds**

| | Skip-gram | Skip-gram topic model |
|---|---|---|
| Naive Bayes | For each word in the corpus $w_i$ | For each word in the corpus $w_i$ |
| | For each word $w_c \in context(i)$<br>Draw $w_c$ via $p(w_c|w_i) \propto exp(v'_{w_c}{}^\mathsf{T} v_{w_i} + b_{w_c})$ | For each word $w_c \in context(i)$<br>Draw $w_c$ via $p(w_c|w_i) = \text{Discrete}(\phi^{(w_i)})$ |
| Mixed membership | For each word in the corpus $w_i$ | For each word in the corpus $w_i$ |
| | Draw a topic $z_i \sim \text{Discrete}(\theta^{(w_i)})$<br>For each word $w_c \in context(i)$<br>Draw $w_c$ via $p(w_c|z_i) \propto exp(v'_{w_c}{}^\mathsf{T} v_{z_i} + b_{w_c})$ | Draw a topic $z_i \sim \text{Discrete}(\theta^{(w_i)})$<br>For each word $w_c \in context(i)$<br>Draw $w_c$ via $p(w_c|z_i) = \text{Discrete}(\phi^{(z_i)})$ |

Table 2: "Generative" models. Identifying the skip-gram (top-left)'s word distributions with topics yields analogous topic models (right), and mixed membership modeling extensions (bottom).

for each document selects a multinomial distribution over words, referred to as a *topic*, with which the documents' words are assumed to be generated. LDA topic models improve over naive Bayes by using a *mixed membership* model, in which the assumption that all words in a document $d$ belong to the same topic is relaxed, and replaced with a *distribution* over topics $\theta^{(d)}$. In the model's assumed generative process, for each word $i$ in document $d$, a topic assignment $z_i$ is drawn via $\theta^{(d)}$, then the word is drawn from the chosen topic $\phi^{(z_i)}$. The mixed membership formalism provides a useful compromise between model flexibility and statistical efficiency: the $K$ topics $\phi^{(k)}$ are shared across all documents, thereby sharing statistical strength, but each document is free to use the topics to its own unique degree. Bayesian inference further aids data efficiency, as uncertainty over $\theta^{(d)}$ can be managed for shorter documents. Some recent papers have aimed to combine topic models and word embeddings (Das et al., 2015; Liu et al., 2015), but they do not aim to address the small data problem for computational social science, which I focus on here. I provide a more detailed discussion of related work in the supplementary.

## 3 The Mixed Membership Skip-Gram

To design an interpretable word embedding model for small corpora, we identify novel connections between word embeddings and topic models, and adapt advances from topic modeling. Following the *distributional hypothesis* (Harris, 1954), the skip-gram's word embeddings parameterize discrete probability distributions over words $p(w_c|w_i)$ which tend to co-occur, and tend to be semantically coherent – a property leveraged by the Gaussian LDA model of Das et al. (2015). This suggests that these discrete distributions can be reinterpreted as *topics* $\phi^{(w_i)}$. We thus reinterpret the skip-gram as a parameterization of a certain supervised naive Bayes topic model (Table 2, top-right). In this topic model, input words $w_i$ are fully observed "cluster assignments," and the words in $w_i$'s contexts are a "document." The skip-gram differs from this supervised

topic model only in the parameterization of the "topics" via word vectors which encode the distributions with a log-bilinear model. Note that although the skip-gram is discriminative, in the sense that it does not jointly model the input words $w_i$, we are here equivalently interpreting it as encoding a "conditionally generative" process for the context given the words, in order to develop probabilistic models that extend the skip-gram.

As in LDA, this model can be improved by replacing the naive Bayes assumption with a mixed membership assumption. By applying the mixed membership representation to this topic model version of the skip-gram, we obtain the model in the bottom-right of Table 2.[1] After once again parameterizing this model with word embeddings, we obtain our final model, the *mixed membership skip-gram (MMSG)* (Table 2, bottom-left). In the model, each input word has a distribution over topics $\theta^{(w)}$. Each topic has a vector-space embedding $v_k$ and each output word has a vector $v'_w$ (a parameter, not an embedding for $w$). A topic $z_i \in \{1, \ldots, K\}$ is drawn for each context, and the words in the context are drawn from the log-bilinear model using $v_{z_i}$:

$$z_i \sim \text{Discrete}(\theta^{(w_i)}) \tag{1}$$

$$p(w_c|z_i) \propto exp(v'_{w_c}{}^\mathsf{T} v_{z_i} + b_{w_c}) . \tag{2}$$

We can expect that the resulting *mixed membership word embeddings* are beneficial in the small-to-medium data regime for the following reasons:

1. By using **fewer input vectors than words**, we can reduce the size of the semantic representation to be learned (output vectors $v'_w$ are viewed as weight parameters, and not used for embedding).

2. The topic vectors are shared across all words, allowing more data to be used per vector.

3. Polysemy is addressed by clustering the words into topics, which leads to topically focused and semantically coherent vector representations.

---

[1] The model retains a naive Bayes assumption at the context level, for latent variable count parsimony.
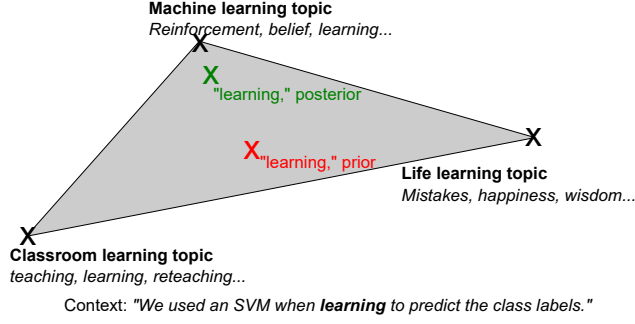
Figure 1: Mixed membership word embeddings $\bar{v}_w$ for word type $w$ (prior) and $\hat{v}_{w_i}$ for word token $w_i$ (posterior), are convex combinations of topic embeddings $v_k$.

Of course, the model also requires some new parameters to be learned, namely the mixed membership proportions $\theta^{(w)}$. Based on topic modeling, I hypothesized that with care, these added parameters need not adversely affect performance in the small-medium data regime, for two reasons: 1) we can use a Bayesian approach to effectively manage uncertainty in them, and to marginalize them out, which prevents them being a bottleneck during training; and 2) at test time, using the posterior for $z_i$ given the context, instead of the "prior" $p(z_i|w_i, \theta) = \theta^{(w_i)}$, mitigates the impact of uncertainty in $\theta^{(w_i)}$ due to limited training data:

$$p(z_i = k|w_i, \text{context}(i), \mathbf{V}, \mathbf{V}', \mathbf{b}, \theta) \quad (3)$$
$$\propto \theta_k^{(w_i)} \prod_{c \in \text{context}(i)} \frac{exp(v_{w_c}'^{\mathsf{T}} v_k + b_{w_c})}{\sum_{j'=1}^{V} exp(v_{j'}'^{\mathsf{T}} v_k + b_{j'})} \; .$$

To obtain a vector for a word type $w$, we can use the prior mean, $\bar{v}_w \triangleq \sum_k v_k \theta_k^{(w)}$. For a word token $w_i$, we can leverage its context via the posterior mean, $\hat{v}_{w_i} \triangleq \sum_k v_k p(z_i = k|w_i, \text{context}(i), \mathbf{V}, \mathbf{V}', \mathbf{b}, \theta)$. These embeddings are convex combinations of topic vectors (see Figure 1 for an example). With fewer vectors than words, some model capacity is lost, but the flexibility of the mixed membership representation allows the model to compensate. When the number of shared vectors equals the number of words, the mixed membership skip-gram is strictly more representationally powerful than the skip-gram. With more vectors than words, we can expect that the increased representational power would be beneficial in the big data regime. As this is not my goal, I leave this for future work.

## 4 Training Algorithm for the MMSG

I first describe an idealized but impractical training algorithm for the MMSG, and then introduce a more practicable procedure (Algorithm 1). The MMSG can in principle be trained via maximum likelihood estimation using EM. Optimizing the log-likelihood is hin-

---

**Algorithm 1** Training the mixed membership skip-gram via annealed MHW and NCE

> **for** $j = 1 : \text{maxAnnealingIter}$ **do**
>    $T_j := T_0 + \lambda \kappa^j$
>    **for** $i = 1 : \text{N}$ **do**
>       $c \sim \text{Uniform}(|\text{context}(w_i)|);$
>       $z_i^{(new)} \sim q_{w_c};$ //via cached alias table samples
>       accept or reject $z_i^{(new)}$ via Equation 6;
>       If accept, $z_i := z_i^{(new)};$
>    **end for**
> **end for**
> $\hat{\theta}_k^{(w_i)} :\propto n_k^{(w_i)\neg i} + \alpha_k$
> $[\mathbf{V}, \mathbf{V}', b] := \text{NCE}(inputWords = \mathbf{z},$
>                 $contextWords = \mathbf{w});$

---

dered by the latent variables, which EM circumvents by focusing on the complete-data log-likelihood (CDLL), $\log p(\mathbf{w}, \mathbf{z}|\mathbf{V}, \mathbf{V}', \mathbf{b}, \theta) =$

$$\sum_{i=1}^{N} \sum_{k=1}^{K} z_{i,k} \log \theta_k^{(w_i)} + \sum_{i=1}^{N} \sum_{k=1}^{K} z_{i,k} \times \quad (4)$$
$$\sum_{c \in \text{context}(i)} \left( v_{w_c}'^{\mathsf{T}} v_k + b_{w_c} - \log \sum_{j'=1}^{D} exp(v_{j'}'^{\mathsf{T}} v_k + b_{j'}) \right) .$$

The E-step computes the *E-step responsibilities* $\gamma_{i,k}$:

$$\gamma_{i,k} = p(z_{i_k} = 1|w_i, \text{context}(i), \{\mathbf{V}, \mathbf{V}', \mathbf{b}, \theta\}^{(old)}) \; ,$$

where $(old)$ superscripts denote current parameter estimates. The M-step optimizes the lower bound on the log-likelihood obtained by substituting $\gamma_{i,k}$ for $z_{i,k}$ in Equation 4. However, this involves a $O(KD)$ complexity for both the E- and M-steps for each token, where $K$ and $D$ are the number of topics/dictionary words, respectively, and even $O(D)$ per token is considered impractical for word embeddings (Mnih and Teh, 2012; Mikolov et al., 2013a). Instead, I propose an approximation to EM that is *sublinear time* in both $K$ and $D$. We first impute $\mathbf{z}$ using a reparameterization technique, thereby reducing the task to standard word embedding. This can be done in sublinear time using the Metropolis-Hastings-Walker (MHW) algorithm. With an oracle $\hat{\mathbf{z}}$ for $\mathbf{z}$, the log-likelihood $\log p(\mathbf{w}|\mathbf{V}, \mathbf{V}', \mathbf{b}, \theta) = \log \sum_{\mathbf{z}} p(\mathbf{w}, \mathbf{z}|\mathbf{V}, \mathbf{V}', \mathbf{b}, \theta)$ simplifies to the CDLL $\log p(\mathbf{w}, \hat{\mathbf{z}}|\mathbf{V}, \mathbf{V}', \mathbf{b}, \theta)$, as in Equation 4. We then efficiently learn the topic and word embeddings via noise-contrastive estimation (NCE). With enough computation NCE exactly optimizes our CDLL objective function, but avoids computing expensive normalization constants and provides an adjustable computational efficiency knob. The details are described below.

| Model | Input word = "Bayesian"<br>Top words in topic for input word. Top 3 topics for word shown for mixed membership models. |
|---|---|
| SGTM | model networks learning neural bayesian data models approach network framework |
| SG | belief learning framework models methods markov function bayesian based inference |
| | |
| MMSGTM | neural bayesian networks mackay computation framework practical learning weigend backpropagation |
| | model models bayesian prior data parameters likelihood priors structure graphical |
| | monte carlo chain markov sampling mcmc method methods model bayesian |
| MMSG | neural networks weigend bayesian data mackay learning computation practical |
| | probability model data models priors algorithm bayesian likelihood set parameters |
| | carlo monte mcmc chain reversible sampling model posterior |

Table 3: SG = skip-gram, TM = topic model, MM = mixed membership.

## 4.1 Imputing the z's

To derive such an algorithm, the key insight is that our MMSG model (Table 2, bottom left) is equivalent to the topic model version (Table 2, bottom right), up to the parameterization. With sufficiently high dimensional embeddings, the log-bilinear model can capture any distribution $p(w_c|z_i)$, and so the maximum likelihood embeddings would encode the exact same word distributions as the MLE topics for the topic model, $\phi^{(z_i)}$. However, the topic model admits a collapsed Gibbs sampler (CGS) that efficiently resolves the cluster assignments, which cause the bottleneck during EM. I therefore propose to reparameterize the MMSG as its corresponding topic model for the purposes of imputing the z's. Then, with the z's fixed to the estimate $\hat{z}$, learning the word and topic vectors corresponds to finding the optimal vectors for encoding the $\phi$'s.

This topic model pre-clustering step is reminiscent of Reisinger and Mooney (2010); Huang et al. (2012); Liu et al. (2015), who apply an off-the-shelf clustering algorithm (or LDA) to initially identify different clusters of contexts, and then apply word embedding algorithms on the cluster assignments. However, our clustering is learned based on the word embedding model itself, and clustering at test time is performed via Bayesian reasoning, in Equation 3, rather than via an ad-hoc method. With Dirichlet priors on the parameters, the collapsed Gibbs update is (derivation in the supplement):

$$p(z_i = k|\cdot) \propto \left( n_k^{(w_i)\neg i} + \alpha_k \right) \qquad (5)$$

$$\times \prod_{c=1}^{|\text{context}(i)|} \frac{n_{w_c}^{(k)\neg i} + \beta_{w_c} + n_{w_c}^{(i,c)}}{n^{(k)\neg i} + \sum_{w'} \beta_{w'} + c - 1} ,$$

where $\alpha$ and $\beta$ are parameter vectors for Dirichlet priors over the topic and word distributions, $n_k^{(w_i)}$ and $n_{w_c}^{(k)\neg i}$ are input and output word/topic counts (excluding the current word), and $n_{w_c}^{(i,c)}$ is the number of occurrences of word $w_c$ before the $c$th word in the $i$th context. We scale this algorithm up to thousands of topics using an adapted version of the recently proposed Metropolis-Hastings-Walker algorithm for high-dimensional topic models, which scales sublinearly in $K$ (Li et al., 2014).

The method uses a data structure called an *alias table*, which allows for amortized O(1) time sampling from discrete distributions. A Metropolis-Hastings update is used to correct for approximating the CGS update with a proposal distribution based on these samples. We can interpret the product over the context, which dominates the collapsed Gibbs update, as a *product of experts* (Hinton, 2002), where each word in the context is an expert which weighs in multiplicatively on the update. In order to approximate this via alias tables, we use proposals which approximate the product of experts with a mixture of experts. We select a word $w_c$ uniformly from the context, and the proposal $q_{w_c}$ draws a candidate topic proportionally to the chosen context word's contribution to the update:

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

We can expect these proposals to bear a resemblance to the target distribution, but to be flatter, which is a property we'd generally like in a proposal distribution. The proposal is implemented efficiently by sampling from the experts via the alias table data structure, in amortized O(1) time, rather than in time linear in the sparsity pattern, as in (Li et al., 2014), since the proposal does not involve the sparse term (which is less important in our case). We perform simulated annealing to optimize over the posterior, which is very natural for Metropolis-Hastings. Interpreting the negative log posterior as the energy function for a Boltzmann distribution at temperature $T_j$ for iteration $j$, this is achieved by raising the model part of the Metropolis-Hastings acceptance ratio to the power of $\frac{1}{T_j}$:

$$z_i^{(new)} \sim q_{w_c} , \ p(\text{accept } z_i^{(new)}|\cdot) =$$

$$\min \left( 1, \left( \frac{p(z_i = z_i^{(new)}|\cdot)}{p(z_i = z_i^{(old)}|\cdot)} \right)^{\frac{1}{T_j}} \frac{q_{w_c}(z_i^{(old)})}{q_{w_c}(z_i^{(new)})} \right) . \quad (6)$$

Annealing also helps with mixing, as the standard Gibbs updates can become nearly deterministic. We use a temperature schedule $T_j = T_0 + \lambda \kappa^j$, where $T_0$ is the target final temperature, $\kappa < 1$, and $\lambda$ controls the initial temperature, and therefore mixing in the early

iterations. In my experiments, I use $T_0 = 0.0001$, $\kappa = 0.99$, and $\lambda = |\text{context}|$. The acceptance probability can be computed in time constant in $K$, and sampling is amortized constant time in $K$, so each iteration is in amortized constant time in $K$. Rao-Blackwellized estimates of the mixed membership proportions are obtained from the final sample as $\hat{\theta}_k^{(w_i)} \propto n_k^{(w_i)\neg i} + \alpha_k$.

## 4.2 Learning the Embeddings

Finally, with the topic assignments $\hat{\mathbf{z}}$ imputed and $\theta$ estimated via the topic model, we must learn the embeddings, which is still an expensive $O(D)$ per context for maximum likelihood estimation, i.e. optimizing

$$\log p(\mathbf{w}, \hat{\mathbf{z}} | \vec{\mathbf{V}}, \mathbf{b}, \theta) = \log p(\mathbf{w} | \hat{\mathbf{z}}, \vec{\mathbf{V}}, \mathbf{b}) + \text{const}, \quad (7)$$

where $\vec{\mathbf{V}}$ is the vector of all word and topic embeddings. This same complexity is also an issue for the standard skip-gram, which Mnih and Teh (2012); Mnih and Kavukcuoglu (2013) have addressed using the noise-contrastive estimation (NCE) algorithm of Gutmann and Hyvärinen (2010, 2012). NCE avoids the expensive normalization step, making the algorithm scale sublinearly in the vocabulary size $D$. The algorithm solves unsupervised learning tasks by transforming them into the supervised learning task of distinguishing the data from randomly sampled noise samples, via logistic regression. Supposing that there are $k$ samples from the noise distribution per word-pair example, the NCE objective function for context $i$ is

$$J^{(i)}(\vec{\mathbf{V}}, \mathbf{b}) \triangleq E_{p_d^{(i)}}[\log \sigma(G(w_c; \vec{\mathbf{V}}, w_i, z_i, \mathbf{b}))]$$
$$- kE_{p_n}[\log(1 - \sigma(G(w_c; \vec{\mathbf{V}}, w_i, z_i, \mathbf{b})))] \quad (8)$$

where $p_d^{(i)}$ is the data distribution for words $w_c$ context $i$, and $G(w_c; \vec{\mathbf{V}}, w_i, z_i, \mathbf{b}) \triangleq \log p(w_c | \vec{\mathbf{V}}, w_i, z_i, \mathbf{b}) - \log p_n(w_c)$ is the difference in log-likelihood between the model and the noise distributions. We learn the embeddings by stochastic gradient ascent on the NCE objective. As the number of noise samples tends to infinity, the method increasingly well approximates maximum likelihood estimation, i.e. the stationary points of Equation 8 converge on those of Equation 7 (Gutmann and Hyvärinen, 2010, 2012).

## 5 Experimental Results

The goals of our experiments were to study the relative merits of big data and domain-specific small data, to validate the proposed methods, and to study their applicability for computational social science research.

## 5.1 Quantitative Experiments

I first measured the effectiveness of the embeddings at the *skip-gram's training task*, predicting context words $w_c$ given input words $w_i$. This task measures the methods' performance for predictive language modeling. I used four datasets of sociopolitical, scientific, and literary interest: the corpus of NIPS articles from 1987 – 1999 ($N \approx 2.3$ million), the U.S. presidential state of the Union addresses from 1790 – 2015 ($N \approx 700,000$), the complete works of Shakespeare ($N \approx 240,000$; this version did not contain the Sonnets), and the writings of black scholar and activist W.E.B. Du Bois, as digitized by Project Gutenberg ($N \approx 170,000$). For each dataset, I held out 10,000 $(w_c, w_i)$ pairs uniformly at random, where $w_c \in \text{context}(i)$, $|\text{context}(i)| = 10$, and aimed to predict $w_c$ given $w_i$ (and optionally, $\text{context}(i) \setminus w_c$). Since there are a large number of classes, I treat this as a ranking problem, and report the mean reciprocal rank. The experiments were repeated and averaged over 5 train/test splits.

The results are shown in Table 4. I compared to a word frequency baseline, the skip-gram (SG), and Tomas Mikolov/Google's vectors trained on Google News, $N \approx 100$ billion, via CBOW. Simulated annealing was performed for 1,000 iterations, NCE was performed for 1 million minibatches of size 128, and 128-dimensional embeddings were used (300 for Google). I used $K = 2,000$ for NIPS, $K = 500$ for state of the Union, and $K = 100$ for the two smaller datasets. Methods were able to leverage the remainder of the context, either by adding the context's vectors, or via the posterior (Equation 3), which helped for all methods except the naive skip-gram. We can identify several noteworthy findings. First, the generic big data vectors (Google+context) were outperformed by the skip-gram on 3 out of 4 datasets (and by the skip-gram topic model on the other), by a large margin, indicating that domain-specific embeddings are often important. Second, the mixed membership models, using posterior inference, beat or matched their naive Bayes counterparts, for both the word embedding models and the topic models. As hypothesized, posterior inference on $z_i$ at test time was important for good performance. Finally, the topic models beat their corresponding word embedding models at prediction. I therefore recommend the use of our MMSG topic model variant for predictive language modeling in the small data regime.

### 5.1.1 Downstream Tasks

I tested the performance of the representations as features for document categorization and regression tasks. The results are given in Table 5. For document categorization, I used three standard benchmark datasets:

| Dataset | Frequency baseline | Google +context | SG | SG +context | MMSG prior | MMSG posterior | SGTM | SGTM +context | MMSGTM prior | MMSGTM posterior |
|---------|-----------|-----------|-------|-------|-------|-------|-------|-------|-------|-------|
| NIPS | 0.029 | 0.027 | 0.038 | 0.031 | 0.037 | **0.062** | **0.055** | **0.064** | **0.046** | **0.074** |
| SOTU | 0.021 | 0.021 | 0.025 | 0.023 | 0.022 | **0.034** | **0.036** | **0.046** | **0.032** | **0.045** |
| Shakespeare | 0.015 | **0.032** | 0.020 | 0.010 | 0.015 | 0.019 | **0.025** | **0.043** | 0.020 | **0.025** |
| Du Bois | 0.028 | 0.033 | 0.045 | 0.037 | 0.041 | **0.053** | **0.052** | **0.081** | **0.050** | **0.066** |

Table 4: Mean reciprocal rank of held-out context words. SG = skip-gram, TM = topic model, MM = mixed membership. Bold indicates statistically significant improvement versus SG.

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

Table 5: Document categorization (top, classification accuracy, larger is better), and predicting the year of State of the Union addresses (bottom, RMSE, LOO cross-validation, smaller is better).

*20 Newsgroups* (19,997 newsgroup posts), *Reuters-150* newswire articles (15,500 articles and 150 classes), and *Ohsumed* medical abstracts on 23 cardiovascular diseases (20,000 articles).[2] I held out 4,000 test documents for 20 Newsgroups, and used the standard train/test splits from the literature in the other corpora (e.g. for *Ohsumed*, 50% of documents were assigned to training and to test sets). I obtained document embeddings for the MMSG, in the same latent space as the topic embeddings, by summing the posterior mean vectors $\hat{v}_{w_i}$ for each token. Vector addition was similarly used to construct document vectors for the other embedding models. All vectors were normalized to unit length. I also considered a tf-idf baseline. Logistic regression models were trained on the features extracted on the training set for each method.

Across the three datasets, several clear trends emerged (Table 5). First, the generic Google vectors were consistently and substantially outperformed in classification performance by the skipgram (SG) and MMSG vectors, highlighting the importance of corpus-specific embeddings. Second, despite the MMSG's superior performance at language modeling on small datasets, the SG features outperformed the MMSG's at the document categorization task. By encoding vectors at the topic level instead of the word level, the MMSG loses word level resolution in the embeddings, which turned out to be valuable for these particular classification tasks. We are not, however, restricted to use only one type of embedding to construct features for classification. Interestingly, when the SG and MMSG features were concatenated (SG+MMSG), this improved classification performance over these vectors individually. This suggests that the topic-level MMSG vectors and word-level SG vectors encode *complementary* information,

and both are beneficial for performance. Finally, further concatenating the generic Google vectors' features (SG+MMSG+Google) improved performance again, despite the fact that these vectors performed poorly on their own. It should be noted that tf-idf, which is notoriously effective for document categorization, outperformed the embedding methods on these datasets.

I also analyzed the regression task of predicting the year of a state of the Union address based on its text information. I used lasso-regularized linear regression models, evaluated via a leave-one-out cross-validation experimental setup. Root-mean-square error (RMSE) results are reported in Table 5 (bottom). Unlike for the other tasks, the Google big data vectors were the best individual features in this case, outperforming the domain-specific SG and MMSG embeddings individually. On the other hand, SG+MMSG+Google performed the best overall, showing that domain-specific embeddings can improve performance even when big data embeddings are successful. The tf-idf baseline was beaten by all of the embedding models on this task.

### 5.2 Computational Social Science Case Studies: State of the Union and NIPS

I also performed several case studies. I obtained document embeddings, in the same latent space as the topic embeddings, by summing the posterior mean vectors $\hat{v}_{w_i}$ for each token, and visualized them in two dimensions using $t$-SNE (Maaten and Hinton, 2008) (all vectors were normalized to unit length). The state of the Union addresses (Figure 2) are embedded almost linearly by year, with a major jump around the New Deal (1930s), and are well separated by party at any given time period. The embedded topics (gray) allow us to interpret the space. The George W. Bush addresses are embedded near a "war on terror" topic ("weapons,

---

[2]All document categorization datasets were obtained from http://disi.unitn.it/moschitti/corpora.htm.
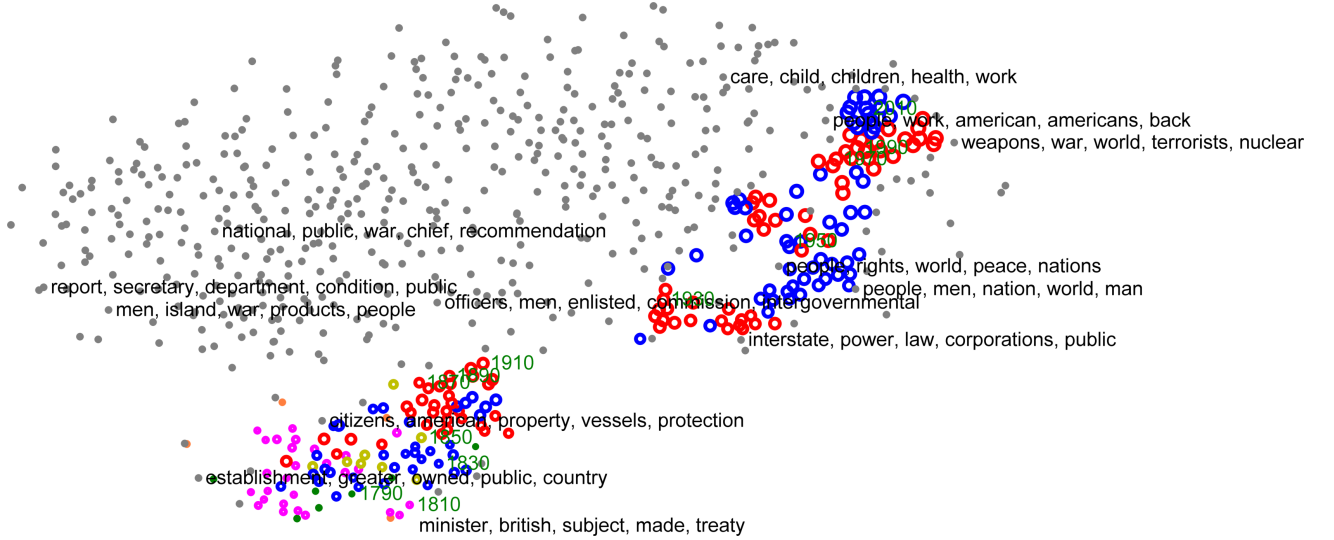
Figure 2: State of the Union (SOTU) addresses. Colored circles are *t*-SNE projected embeddings for SOTU addresses. Color = party (red = GOP, blue = Democrats, light green = Whigs, pink = Democratic-Republicans, orange = Federalists (John Adams), green = George Washington), size = recency (year, see dates in green). Gray circles correspond to topics.

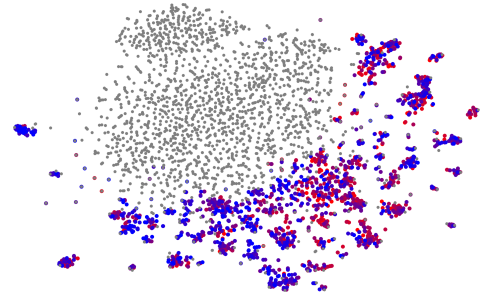| Nearest topic after composition of mean vectors for words | |
|---|---|
| object + recognition | objects visual object recognition model |
| character + recognition | recognition segmentation character |
| speech + recognition | speech recognition hmm system hybrid |
| computer + vision | computer vision ieee image pattern |
| computer + science | university science colorado department |
| bias + variance | error training set data performance |
| covariance + variance | gaussian distribution model matrix |



Figure 3: **Left:** Vector compositionality examples, NIPS. **Right:** NIPS documents/ topics, *t*-SNE.

war..."), and the Barack Obama addresses are embedded near a "stimulus" topic ("people, work...").

On the NIPS corpus, for the input word "Bayesian" (Table 3), the naive Bayes and skip-gram models learned a topic with words that refer to Bayesian networks, probabilistic models, and neural networks. The mixed membership models are able to separate this into more coherent and specific topics including Bayesian modeling, Bayesian training of neural networks (for which Sir David MacKay was a strong proponent, and Andreas Weigend wrote an influential early paper), and Monte Carlo methods. By performing the additive composition of word vectors, which we obtain by finding the prior mean vector for each word type $w$, $\bar{v}_w \triangleq \sum_k v_k \theta_k^{(w)}$ (and then normalizing), we obtain relevant topics $v_k$ as nearest neighbors (Figure 3). Similarly, we find that the additive composition of topic and word vectors works correctly: $v_{objectRecognition} - \bar{v}_{object} + \bar{v}_{speech} \approx v_{speechRecognition}$, and $v_{speechRecognition} - \bar{v}_{speech} + \bar{v}_{character} \approx v_{characterRecognition}$.

The *t*-SNE visualization of NIPS documents (Figure 3) shows some temporal clustering patterns (blue documents are more recent, red documents are older, and gray points are topics). I provide a more detailed case study on NIPS in the supplementary material.

# 6 Conclusion

I have proposed a model-based method for training interpretable corpus-specific word embeddings for computational social science, using mixed membership representations, Metropolis-Hastings-Walker sampling, and NCE. Experimental results for prediction, supervised learning, and case studies on state of the Union addresses and NIPS articles, indicate that high-quality embeddings and topics can be obtained using the method. The results highlight the fact that big data is not always best, as domain-specific data can be very valuable, even when it is small. I plan to use this approach for substantive social science applications, and to address algorithmic bias and fairness issues.

# References

Airoldi, E., Blei, D., Feinberg, S., and Xing, E. (2008). Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9:1981–2014.

Airoldi, E. M., Blei, D. M., Erosheva, E. A., and Fienberg, S. E. (2014). Introduction to mixed membership models and methods. In *Handbook of Mixed Membership Models and Their Applications*. Chapman and Hall/CRC.

Barnard, K., Duygulu, P., Forsyth, D., Freitas, N. d., Blei, D. M., and Jordan, M. I. (2003). Matching words and pictures. *Journal of Machine Learning Research*, 3(Feb):1107–1135.

Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language model. *Journal of Machine Learning Research*, 3:1137–1155.

Blei, D., Ng, A., and Jordan, M. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022.

Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. In *Advances in Neural Information Processing Systems 29*, pages 4349–4357.

Brent, M. R. (1999). An efficient, probabilistically sound algorithm for segmentation and word discovery. *Machine Learning*, 34(1):71–105.

Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537.

Das, R., Zaheer, M., and Dyer, C. (2015). Gaussian LDA for topic models with word embeddings. In *Proceedings of the 53nd Annual Meeting of the Association for Computational Linguistics*, pages 795–804.

Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6):391.

Erosheva, E., Fienberg, S., and Lafferty, J. (2004). Mixed-membership models of scientific publications. *Proceedings of the National Academy of Sciences of the United States of America*, 101(Suppl 1):5220–5227.

Erosheva, E. A. (2003). Bayesian estimation of the grade of membership model. *Bayesian Statistics*, 7:501–510.

Fei-Fei, L. and Perona, P. (2005). A Bayesian hierarchical model for learning natural scene categories. In *Proceedings of the 2005 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 524–531. IEEE.

Griffiths, T. L., Steyvers, M., and Tenenbaum, J. B. (2007). Topics in semantic representation. *Psychological Review*, 114(2):211.

Grimmer, J. (2010). A Bayesian hierarchical topic model for political texts: Measuring expressed agendas in senate press releases. *Political Analysis*, pages 1–35.

Guo, F., Blundell, C., Wallach, H., and Heller, K. (2015). The Bayesian echo chamber: Modeling social influence via linguistic accommodation. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 315–323.

Gutmann, M. and Hyvärinen, A. (2010). Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*.

Gutmann, M. U. and Hyvärinen, A. (2012). Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *Journal of Machine Learning Research*, 13(Feb):307–361.

Harris, Z. S. (1954). Distributional structure. *Word*, 10(2-3):146–162.

Hinton, G. E. (2002). Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800.

Hinton, G. E., Mcclelland, J. L., and Rumelhart, D. E. (1986). Distributed representations. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Volume 1: Foundations*, chapter 3, pages 77–109. MIT Press, Cambridge, MA.

Hinton, G. E. and Salakhutdinov, R. R. (2009). Replicated softmax: an undirected topic model. In *Advances in Neural Information Processing Systems*, pages 1607–1614.

Hofmann, T. (1999a). Probabilistic latent semantic analysis. In *Proceedings of the Fifteenth International Conference on Uncertainty in Artificial Intelligence*, pages 289–296. Morgan Kaufmann Publishers Inc.

Hofmann, T. (1999b). Probabilistic latent semantic indexing. In *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 50–57. ACM.

Huang, E. H., Socher, R., Manning, C. D., and Ng, A. Y. (2012). Improving word representations via global context and multiple word prototypes. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 873–882. Association for Computational Linguistics.

Kiros, R., Zhu, Y., Salakhutdinov, R. R., Zemel, R., Urtasun, R., Torralba, A., and Fidler, S. (2015). Skip-thought vectors. In *Advances in Neural Information Processing Systems*, pages 3294–3302.

Larochelle, H. and Lauly, S. (2012). A neural autoregressive topic model. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 25*, pages 2708–2716.

Li, A. Q., Ahmed, A., Ravi, S., and Smola, A. J. (2014). Reducing the sampling complexity of topic models. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 891–900. ACM.

Liu, Y., Liu, Z., Chua, T.-S., and Sun, M. (2015). Topical word embeddings. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2418–2424.

Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(Nov):2579–2605.

Manton, K. G., Tolley, H. D., and Woodbury, M. A. (1994). *Statistical applications using fuzzy sets*. Wiley-Interscience.

Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. *Proceedings of the 2013 International Conference on Learning Representations (ICLR)*.

Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119.

Mimno, D. (2012). Computational historiography: Data mining in a century of classics journals. *Journal on Computing and Cultural Heritage (JOCCH)*, 5(1):3.

Mnih, A. and Kavukcuoglu, K. (2013). Learning word embeddings efficiently with noise-contrastive estimation. In *Advances in Neural Information Processing Systems*, pages 2265–2273.

Mnih, A. and Teh, Y. W. (2012). A fast and simple algorithm for training neural probabilistic language models. In *Proceedings of the 29th International Conference on Machine Learning*.

Nguyen, V.-A., Boyd-Graber, J., Resnik, P., Cai, D. A., Midberry, J. E., and Wang, Y. (2014). Modeling topic control to detect influence in conversations using nonparametric topic models. *Machine Learning*, 95(3):381–421.

Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, volume 14, pages 1532–1543.

Pritchard, J. K., Stephens, M., and Donnelly, P. (2000). Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959.

Reisinger, J. and Mooney, R. J. (2010). Multi-prototype vector-space models of word meaning. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 109–117. Association for Computational Linguistics.

Roberts, M. E., Stewart, B. M., Tingley, D., Lucas, C., Leder-Luis, J., Gadarian, S. K., Albertson, B., and Rand, D. G. (2014). Structural topic models for open-ended survey responses. *American Journal of Political Science*, 58(4):1064–1082.

Teh, Y. W., Jordan, M. I., Beal, M. J., and Blei, D. M. (2006). Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476).

Tian, F., Dai, H., Bian, J., Gao, B., Zhang, R., Chen, E., and Liu, T.-Y. (2014). A probabilistic model for learning multi-prototype word embeddings. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, pages 151–160.

Vaswani, A., Zhao, Y., Fossum, V., and Chiang, D. (2013). Decoding with large-scale neural language models improves translation. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1387–1392.

Wallach, H. (2016). Computational social science: Toward a collaborative future. In Alvarez, R. M., editor, *Computational Social Science: Discovery and Prediction*. Cambridge University Press.

Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., and Fidler, S. (2015). Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 19–27.

# Supplementary Material

## A    Related Work

In this supplementary document, we discuss related work in the literature and its relation to our proposed methods, provide a case study on NIPS articles, and derive the collapsed Gibbs sampling update for the MMSGTM, which we leverage when training the MMSG.

### A.1    Topic Modeling and Word Embeddings

The *Gaussian LDA* model of Das et al. (2015) improves the performance of topic modeling by leveraging the semantic information encoded in word embeddings. Gaussian LDA modifies the generative process of LDA such that each topic is assumed to generate the vectors via its own Gaussian distribution. Similarly to our MMSG model, in Gaussian LDA each topic is encoded with a vector, in this case the mean of the Gaussian. It takes pre-trained word embeddings as input, rather than learning the embeddings from data within the same model, and does not aim to perform word embedding.

The topical word embedding (TWE) models of Liu et al. (2015) reverse this, as they take LDA topic assignments of words as input, and aim to use them to improve the resultant word embeddings. The authors propose three variants, each of which modifies the skip-gram training objective to use LDA topic assignments together with words. In the best performing variant, called *TWE-1*, a standard skip-gram word embedding model is trained independently with another skip-gram variant, which tries to predict context words given the input word's topic assignment. The skip-gram embedding and the topic embeddings are concatenated to form the final embedding.

At test time, a distribution over topics for the word given the context, $p(z_i|\text{context}(i))$ is estimated according to the topic counts over the other context words. Using this as a prior, a posterior over topics given both the input word and the context is calculated, and similarities between pairs of words (with their contexts) are averaged over this posterior, in a procedure inspired by those used by Reisinger and Mooney (2010); Huang et al. (2012). The primary similarity to our MMSG approach is the use of a training algorithm involving the prediction of context words, given a topic. Our method does this as part of an overall model-based inference procedure, and we learn mixed membership proportions $\theta^{(w)}$ rather than using empirical counts as the prior over topics for a word token. In accordance with the skip-gram's prediction model, we are thus able to model the context words in the data likelihood term when computing the posterior probability of the topic assignment. TWE-1 requires that topic assignments are available at test time. It provides a mechanism to predict contextual similarity, but not to predict held-out context words, so we are unable to compare to it in our experiments.

Other neurally-inspired topic models include replicated softmax (Hinton and Salakhutdinov, 2009), and its successor, DocNADE (Larochelle and Lauly, 2012). Replicated softmax extends the restricted Boltzmann machine to handle multinomial counts for document modeling. DocNADE builds on the ideas of replicated softmax, but uses the NADE architecture, where observations (i.e. words) are modeled sequentially given the previous observations.

### A.2    Multi-Prototype Embedding Models

Multi-prototype embeddings models are another relevant line of work. These models address lexical ambiguity by assigning multiple vectors to each word type, each corresponding to a different meaning of that word. Reisinger and Mooney (2010) propose to cluster the occurrences of each word type, based on features extracted from its context. Embeddings are then learned for each cluster. Huang et al. (2012) apply a similar approach, but they use initial single-prototype word embeddings to provide the features used for clustering. These clustering methods have some resemblance to our topic model pre-clustering step, although their clustering is applied within instances of a given word type, rather than globally across all word types, as in our methods. This results in models with more vectors than words, while we aim to find fewer vectors than words, to reduce the model's complexity for small datasets. Rather than employing an off-the-shelf clustering algorithm and then applying an unrelated embedding model to its output, our approach aims to perform model-based clustering within an overall joint model of topic/cluster assignments and word vectors.

Perhaps the most similar model to ours in the literature is the probabilistic multi-prototype embedding model of Tian et al. (2014), who treat the prototype assignment of a word as a latent variable, assumed drawn from a mixture over prototypes for each word. The embeddings are then trained using EM. Our MMSG model can be understood as the mixed membership version of this model, in which the prototypes (vectors) are shared across all word types, and each word type has its own mixed membership proportions across the shared prototypes. While a similar EM algorithm can be applied to the MMSG, the E-step is much more expensive, as we typically desire many more shared vectors (often in the thousands) than we would prototypes per a single word type (Tian et al. use ten in their experiments). We use the Metropolis-Hastings-Walker algorithm with the topic model reparameterization of our model in order to address this by efficiently pre-solving the E-step.

### A.3    Mixed Membership Modeling

Mixed membership modeling is a flexible alternative to traditional clustering, in which each data point is assigned to a single cluster. Instead, mixed membership models posit that individual entities are associated with multiple underlying clusters, to differing degrees, as encoded by a mixed membership vector that sums to one across the clusters (Erosheva et al., 2004; Airoldi et al., 2014). These mixed membership proportions are generally used to model lower-level grouped data, such as the words inside a document. Each lower-level data point inside a group is assumed to be assigned to one of the shared, global clusters according to the group-level membership proportions. Thus, a mixed membership model consists of a mixture model for each group, which share common mixture component parameters,

Figure 4: NIPS documents/topics, $t$-SNE, zoomed in. Blue/red = more recent/older, gray = topics.



Figure 5: NIPS authors and topics, $t$-SNE, zoomed in. Blue = authors, gray = topics.

but with differing mixture proportions.

This formalism has lead to probabilistic models for a variety of applications, including medical diagnosis (Manton et al., 1994), population genetics (Pritchard et al., 2000), survey analysis (Erosheva, 2003), computer vision (Barnard et al., 2003; Fei-Fei and Perona, 2005), text documents (Hofmann, 1999b; Blei et al., 2003), and social network analysis (Airoldi et al., 2008). Nonparametric Bayesian extensions, in which the number of underlying clusters is learned from data via Bayesian inference, have also been proposed (Teh et al., 2006). In this work, dictionary words are assigned a mixed membership distribution over a set of shared latent vector space embeddings. Each instantiation of a dictionary word (an "input" word) is assigned to one of the shared embeddings based on its dictionary word's membership vector. The words in its context ("output" words) are assumed to be drawn based on the chosen embedding.

## B  Case Study on NIPS

In Figure 4, we show a zoomed in $t$-SNE visualization of NIPS document embeddings. We can see regions of the space corresponding to learning algorithms (bottom), data space and latent space (center), training neural networks (top), and nearest neighbors (bottom-left). We also visualized the authors' embeddings via $t$-SNE (Figure 5). We find regions of latent space for reinforcement learning authors (left: "state, action,...," Singh, Barto,Sutton), probabilistic methods (right: "mixture, model," "monte, carlo," Bishop, Williams, Barber, Opper, Jordan, Ghahramani, Tresp, Smyth), and evaluation (top-right: "results, performance, experiments,...").

## C  Derivation of the Collapsed Gibbs Update

Let $C_i = |\text{context}(i)|$ be the number of output words in the $i$th context, let $w_1^{(i)}, \ldots, w_{C_i}^{(i)}$ be those output words,

and let $\mathbf{w}_{\neg i}$ be the input words other that $w_i$ (similarly, topic assignments $\mathbf{z}_{\neg i}$ and output words $\mathbf{w}^{(\neg i)}$). Then the collapsed Gibbs update samples from the conditional distribution

$$
\begin{aligned}
& p(z_i = k | \mathbf{z}_{\neg i}, w_i, w_1^{(i)}, \ldots, w_{C_i}^{(i)}, \mathbf{w}_{\neg i}, \mathbf{w}^{(\neg i)}, \alpha, \beta) \\
& \propto p(z_i = k, w_1^{(i)}, \ldots, w_{C_i}^{(i)} | \mathbf{z}_{\neg i}, w_i, \mathbf{w}_{\neg i}, \mathbf{w}^{(\neg i)}, \alpha, \beta) \\
& = \int_{\phi^{(k)}} \int_{\theta^{(w_i)}} p(z_i = k, w_1^{(i)}, \ldots, w_{C_i}^{(i)}, \phi^{(k)}, \theta^{(w_i)} | \mathbf{z}_{\neg i}, \\
& \qquad\qquad\qquad w_i, \mathbf{w}_{\neg i}, \mathbf{w}^{(\neg i)}, \alpha, \beta) \\
& = \int_{\phi^{(k)}} \int_{\theta^{(w_i)}} p(z_i = k, w_1^{(i)}, \ldots, w_{C_i}^{(i)} | \phi^{(k)}, \theta^{(w_i)}, w_i) \\
& \qquad\qquad \times p(\phi^{(k)}, \theta^{(w_i)} | \mathbf{z}_{\neg i}, w_i, \mathbf{w}_{\neg i}, \mathbf{w}^{(\neg i)}, \alpha, \beta) \\
& = \int_{\phi^{(k)}} \int_{\theta^{(w_i)}} \theta_k^{(w_i)} \prod_{c=1}^{C_i} \phi_{w_c^{(i)}}^{(k)} \times p(\theta^{(w_i)} | \mathbf{z}_{\neg i: w_j = w_i}, \alpha) \\
& \qquad\qquad \times p(\phi^{(k)} | \mathbf{z}_{\neg i}, \mathbf{w}^{(\neg i)}, \beta) \\
& = \int_{\theta^{(w_i)}} \theta_k^{(w_i)} p(\theta^{(w_i)} | \mathbf{z}_{\neg i: w_j = w_i}, \alpha) \\
& \qquad \times \int_{\phi^{(k)}} \prod_{c=1}^{C_i} \phi_{w_c^{(i)}}^{(k)} p(\phi^{(k)} | \mathbf{z}_{\neg i}, \mathbf{w}^{(\neg i)}, \beta) \; .
\end{aligned}
$$

We recognize the first integral as the mean of a Dirichlet distribution which we obtain via conjugacy:

$$
p(\theta^{(w_i)} | \mathbf{z}_{\neg i: w_j = w_i}, \alpha) = \text{Dirichlet}(\mathbf{n}_{\cdot}^{(w_i) \neg i} + \alpha)
$$

$$
\begin{aligned}
\int_{\theta^{(w_i)}} \theta_k^{(w_i)} p(\theta^{(w_i)} | \mathbf{z}_{\neg i: w_j = w_i}, \alpha) & = \frac{n_k^{(w_i) \neg i} + \alpha_k}{\sum_{k'} n_{k'}^{(w_i) \neg i} + \alpha_{k'}} \\
& \propto n_k^{(w_i) \neg i} + \alpha_k \; .
\end{aligned}
$$

The above can also be understood as the probability of the next ball drawn from a multivariate Polya urn model, also known as the Dirichlet-compound multinomial distribution, arising from the posterior predictive distribution of a discrete likelihood with a Dirichlet prior. We will need the full form of such a distribution to analyze the second integral. Once again leveraging conjugacy, we have:

$$
\int_{\phi^{(k)}} \prod_{c=1}^{C_i} \phi_{w_c^{(i)}}^{(k)} p(\phi^{(k)} | \mathbf{z}_{\neg i}, \mathbf{w}^{(\neg i)}, \beta)
$$

$$
= \int_{\phi^{(k)}} \prod_{c=1}^{C_i} \phi_{w_c^{(i)}}^{(k)} \frac{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v))}{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v)} \prod_{v=1}^D \phi_v^{(k) n_v^{(k) \neg i} + \beta_v - 1}
$$

$$
= \int_{\phi^{(k)}} \frac{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v))}{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v)} \prod_{v=1}^D \phi_v^{(k) n_v^{(k) \neg i} + \beta_v + n_v^{(i)} - 1}
$$

$$
\begin{aligned}
& = \frac{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v))}{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v)} \frac{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v + n_v^{(i)})}{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v + n_v^{(i)}))} \\
& \times \int_{\phi^{(k)}} \frac{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v + n_v^{(i)}))}{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v + n_v^{(i)})} \prod_{v=1}^D \phi_v^{(k) n_v^{(k) \neg i} + \beta_v + n_v^{(i)} - 1} \\
& = \frac{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v))}{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v)} \frac{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v + n_v^{(i)})}{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v + n_v^{(i)}))} \; ,
\end{aligned}
$$

where $n_v^{(i)}$ is the number of times that output word $v$ occurs in the $i$th context, since the final integral is over the full support of a Dirichlet distribution, which integrates to one. Eliminating terms that aren't affected by the $z_i$ assignment, the above is

$$
\begin{aligned}
& \propto \frac{\prod_{v=1}^D \Gamma(n_v^{(k) \neg i} + \beta_v + n_v^{(i)})}{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v + n_v^{(i)}))} \\
& = \frac{\prod_{v=1}^D \left( \Gamma(n_v^{(k) \neg i} + \beta_v) \prod_{j=0}^{n_v^{(i)} - 1} (n_v^{(k) \neg i} + \beta_v + j) \right)}{\Gamma(\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v)) \prod_{j=0}^{C_i - 1} (\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v) + j)} \\
& \propto \frac{\prod_{v=1}^D \prod_{j=0}^{n_v^{(i)} - 1} (n_v^{(k) \neg i} + \beta_v + j)}{\prod_{j=0}^{C_i - 1} (\sum_{v=1}^D (n_v^{(k) \neg i} + \beta_v) + j)} \\
& = \prod_{c=1}^{C_i} \frac{n_{w_c}^{(k) \neg i} + \beta_{w_c} + n_{w_c^{(i,c)}}}{n^{(k) \neg i} + \sum_v \beta_v + c - 1}
\end{aligned}
$$

where we have used the fact that $\Gamma(x + n) = (x + n - 1)(x + n - 2)...(x + 1)x\Gamma(x)$ for any $x > 0$, and integer $n \geq 1$. We can interpret this as the probability of drawing the context words under the multivariate Polya urn model, in which the number of "colored balls" (word counts plus prior counts) is increased by one each time a certain color (word) is selected. In other words, in each step, corresponding to the selection of each context word, we draw a ball from the urn, then put it back, *along with another ball of the same color*. The $n_{w_c^{(i,c)}}$ and $c - 1$ terms reflect that the counts have been changed by adding these extra balls into the urn in each step. The second to last equation shows that this process is exchangeable: it does not matter which order the balls were drawn in when determining the probability of the sequence. Multiplying this with the term from the first integral, calculated earlier, gives us the final form of the update equation,

$$
p(z_i = k | \cdot) \propto (n_k^{(w_i) \neg i} + \alpha_k) \prod_{c=1}^{C_i} \frac{n_{w_c}^{(k) \neg i} + \beta_{w_c} + n_{w_j^{(i,c)}}}{n^{(k) \neg i} + \sum_v \beta_v + c - 1} \; .
$$