

This is a pre-print of an article published in Knowledge and Information Systems. The final authenticated version is available online at: <https://doi.org/10.1007/s10115-018-1321-8>
Received: 08 December 2017 / Accepted: 22 November 2018

Learning Concept Embeddings for Dataless Classification via Efficient Bag of Concepts Densification

Walid Shalaby

Wlodek Zadrozny

Computer Science Department,

University of North Carolina at Charlotte,

Charlotte, NC, 28223 USA

WSHALABY@UNCC.EDU

WZADROZN@UNCC.EDU

Abstract

Explicit concept space models have proven efficacy for text representation in many natural language and text mining applications. The idea is to embed textual structures into a semantic space of concepts which captures the main ideas, objects, and the characteristics of these structures. The so called Bag of Concepts (BoC) representation suffers from data sparsity causing low similarity scores between similar texts due to low concept overlap. To address this problem, we propose two neural embedding models to learn continuous concept vectors. Once they are learned, we propose an efficient vector aggregation method to generate fully continuous BoC representations. We evaluate our concept embedding models on three tasks: 1) measuring entity semantic relatedness and ranking where we achieve 1.6% improvement in correlation scores, 2) dataless concept categorization where we achieve state-of-the-art performance and reduce the categorization error rate by more than 5% compared to five prior word and entity embedding models, and 3) dataless document classification where our models outperform the sparse BoC representations. In addition, by exploiting our efficient linear time vector aggregation method, we achieve better accuracy scores with much less concept dimensions compared to previous BoC densification methods which operate in polynomial time and require hundreds of dimensions in the BoC representation.

1. Introduction

Vector space representation models are used to represent textual structures (words, phrases, and documents) as multidimensional vectors. Typically, those models utilize textual corpora and/or Knowledge Bases (KBs) in order to extract and model real-world knowledge. Once acquired, any given text is represented as a *vector* in the semantic space. The goal is thus to accurately place similar structures close to each other in that semantic space, while placing dissimilar ones far apart.

Explicit concept space models are one of these *vector-based representations* which are motivated by the idea that, high level cognitive tasks such learning and reasoning are supported by the knowledge we acquire from *concepts*¹ (Song et al., 2015). Therefore, such models utilize concept vectors (aka bag-of-concepts (BoC)) as the underlying semantic rep-

1. A concept is an expression that denotes an idea, event, or an object.

Category	Top 3 Concepts	Instance Top 3 Concepts	ESA	CCX	CRX
Hockey	<ul style="list-style-type: none"> - Detroit Red Wings, - History of the Detroit Red Wings, - History of the NHL on United States television 	Instance (53798) <ul style="list-style-type: none"> - History of the Detroit Red Wings, - Detroit Red Wings, - Pittsburgh Penguins 	0.73	0.95	0.95
		Instance (54551) <ul style="list-style-type: none"> - Paul Kariya, - Boston Bruins, - Bobby Orr 	0.0	0.84	0.80
Guns	<ul style="list-style-type: none"> - Waco siege, - Overview of gun laws by nation, - Gun violence in the United States 	Instance (54387) <ul style="list-style-type: none"> - Overview of gun laws by nation, - Waco siege, - Gun politics in the United States 	0.71	0.94	0.93
		Instance (54477) <ul style="list-style-type: none"> - Concealed carry in the United States, - Overview of gun laws by nation, - Gun laws in California 	0.33	0.80	0.75

Table 1: Top 3 concepts generated using ESA (Gabrilovich & Markovitch, 2007) for two 20-newsgroups categories (Hockey and Guns) along with top 3 concepts of sample instances. Using exact match similarity scoring (as in ESA) result in low scores between similar instance and category concept vectors. When using concept embeddings (our models), we obtain relatively higher and more representative similarities.

resentation of a given text through a process called *conceptualization*, which is mapping the text into relevant concepts capturing its main ideas, objects, events, and their characteristics. The concept space typically include concepts obtained from KBs such as Wikipedia, Probase (Wu et al., 2012), and others. Once the concept vectors are generated, similarity between two concept vectors can be computed using a suitable similarity/distance measure such as *cosine*.

Similar to the traditional Bag-of-Words (BoW) representation, the BoC vector is a multidimensional *sparse* vector whose dimensionality is the same as the number of concepts in the employed KB (typically *millions*). Consequently, it suffers from *data sparsity* causing low similarity scores between similar texts due to low concept overlap. Formally, given a text snippet $T = \{t_1, t_2, \dots, t_n\}$ of n terms where $n \geq 1$, and a concept space $C = \{c_1, c_2, \dots, c_N\}$ of size N . The BoC vector $\mathbf{v} = \{w_1, w_2, \dots, w_N\} \in \mathbb{R}^N$ of T is a vector of weights of each concept where each w_i of concept c_i is calculated as in equation 1:

$$w_i = \sum_{j=1}^n f(c_i, t_j), 1 \leq i \leq N \quad (1)$$

Here $f(c, t)$ is a *scoring function* which indicates the degree of *association* between term t and concept c . For example, Gabrilovich and Markovitch (2007) proposed Explicit Semantic

Analysis (ESA) which uses Wikipedia articles as concepts and the TF-IDF score of the terms in these article as the association score. Another scoring function might be the co-occurrence count or Pearson correlation score between t and c . As we can notice, only very small subset of the concept space would have nonzero scores with a given term². Moreover, the BoC vector is generated from the top n concepts which have relatively high association scores with the input terms (typically few hundreds). Thus each text snippet is mapped to a very sparse vector of millions of dimensions having only few hundreds nonzero values leading to the *BoC sparsity* problem (Peng et al., 2016).

Typically, the *cosine* similarity measure is used compute the similarity between a pair of BoC vectors \mathbf{u} and \mathbf{v} . Because the concept vectors are very sparse and for space efficiency, we can rewrite each vector as a vector of tuples (c_i, w_i) . Suppose that $\mathbf{u} = \{(c_{n_1}, u_1), \dots, (c_{n_{|\mathbf{u}|}}, u_{|\mathbf{u}|})\}$ and $\mathbf{v} = \{(c_{m_1}, v_1), \dots, (c_{m_{|\mathbf{v}|}}, v_{|\mathbf{v}|})\}$, where u_i and v_j are the corresponding weights of concepts c_{n_i} and c_{m_j} respectively. And n_i, m_j are the indices of these concepts in the concept space C such that $1 \leq n_i, m_j \leq N$. Then, the similarity score can be written as in equation 2:

$$Sim_{cos}(\mathbf{u}, \mathbf{v}) = \frac{\sum_{i=1}^{|\mathbf{u}|} \sum_{j=1}^{|\mathbf{v}|} \mathbb{1}(n_i=m_j) u_i v_j}{\sqrt{\sum_{i=1}^{|\mathbf{u}|} u_i^2} \sqrt{\sum_{j=1}^{|\mathbf{v}|} v_j^2}} \quad (2)$$

where $\mathbb{1}$ is the indicator function which returns 1 if $n_i=m_j$ and 0 otherwise. Having such sparse representation and using exact match similarity scoring measure, we can expect that two very similar text snippets might have *zero similarity* score if they map to *different but very related set of concepts* (Song & Roth, 2015). We demonstrate this fact in Table 1 (ESA column).

In this paper we utilize *neural-based representations* to overcome the BoC sparsity problem. The basic idea is to *map* each concept to a *fixed size continuous vector*³. These vectors can then be used to compute concept-concept similarity and thus overcome the concept mismatch problem.

Our work is also motivated by the success of recent neural-based methods for learning word embeddings in capturing both syntactic and semantic regularities using simple vector arithmetic (Mikolov, Chen, Corrado, & Dean, 2013a; Mikolov, Sutskever, Chen, Corrado, & Dean, 2013b; Pennington, Socher, & Manning, 2014). For example, inferring analogical relationships between words: $vec(king) - vec(man) + vec(woman) = vec(queen)$. This indicates that the learned vectors encode meaningful multi-clustering for each word.

However, word vectors suffer from significant limitations. First, each word is assumed to have a *single meaning* regardless of its context and thus is represented by a single vector in the semantic space (e.g., *charlotte (city)* vs. *charlotte (given name)*). Second, the space contains vectors of single words only. Vectors of multiword expressions (MWEs) are typically obtained by averaging the vectors of individual words. This often produces inaccurate representations especially if the meaning of the MWE is different from the composition of meanings of its individual words (e.g., $vec(north\ carolina)$ vs. $vec(north) + vec(carolina)$). Additionally, mentions that are used to refer to the same concept would have different

2. Unless the term is very common (e.g., a, the, some...etc) and carry no relevant information.

3. We use the terms continuous, dense, distributed vectors interchangeably to refer to real-valued vectors.

embeddings (e.g., *u.s.*, *america*, *usa*), and the model might not be able to place those individual vectors in the same sub-cluster, especially the rare surface forms.

We propose two *neural embedding* models in order to learn continuous concept vectors based on the skip-gram model (Mikolov et al., 2013b). Our first model is the *Concept Raw Context* model (CRX) which utilizes raw concept mentions in a large scale textual KB to jointly learn embeddings of both words and concepts. Our second model is the *Concept-Concept Context* model (CCX) which learns the embeddings of concepts from their conceptual contexts (i.e., contexts containing surrounding concepts only).

After learning the concept vectors, we propose an *efficient BoC aggregation* method. We perform *weighted average* of the individual concept vectors to generate fully *continuous* BoC representations (CBoC). This aggregation method allows measuring the similarity between pairs of BoC in *linear time* which is more efficient than previous methods that require *quadratic* or at least *log-linear* time if optimized (see equation 2). Our embedding models produce more *representative* similarity scores for BoC containing *different but semantically similar* concepts as shown in Table 1 (columns 2-3).

We evaluate our embedding models on three tasks:

1. An intrinsic task of measuring entity semantic relatedness and ranking where we achieve 1.6% improvement in correlation scores.
2. Dataless concept categorization where we achieve state-of-the-art performance and reduce the categorization error rate by more than 5% compared to five prior word and entity embedding models.
3. An extrinsic task of dataless document classification which is a learning *protocol* used to perform text categorization without the need for labeled data to train a classifier (Chang et al., 2008). Experimental results show that we can achieve better accuracy using our efficient BoC densification method compared to the original sparse BoC representation with much less concept dimensions.

The contributions of this paper are fourfold: First, we propose two *low cost* concept embedding models which learn concept representations from concept mentions in free-text corpora. Our models require few hours rather than days to train. Second, we show through empirical results the efficacy of the learned concept embeddings in measuring entity semantic relatedness and concept categorization. Our models achieve *state-of-the-art* performance on two concept categorization datasets. Third, we propose simple and efficient vector aggregation method to obtain *fully dense BoC in linear time*. Fourth, we demonstrate through experiments on dataless document classification that we can obtain better accuracy using the dense BoC representation with much less dimensions (few in most cases), reducing the *computational cost* of generating the BoC vector significantly.

The rest of this paper is organized as follows: Section 2 reviews related work; Section 3 describes our proposed embedding models; Section 4 introduces the applications of the proposed models; Section 5 reports experiments and results; and finally, Section 6 provides discussion and concluding remarks.

2. Related Work

Text Conceptualization: Humans understand languages through multi-step cognitive processes which involves building rich models of the world and making multi-level gen-

eralizations from the input text (Wang & Wang, 2016). One way of automating such generalizations is through text conceptualization. Either by extracting basic level concepts from the input text using concept KBs (Kim et al., 2013; Song et al., 2015), or mapping the whole input into a concept space that capture its semantics as in ESA (Gabrilovich & Markovitch, 2007) and MSA (Shalaby & Zadrozny, 2015).

One major line of conceptualization research utilizes semi-structured KBs such as Wikipedia in order to construct the concept space which is defined by all Wikipedia article titles. Such models have proven efficacy for semantic analysis of textual data especially short texts where contextual information is missing or insufficient (see Shalaby and Zadrozny (2015) for examples).

Another research direction uses more structured concept KBs such as Probase⁴ (Wu et al., 2012). Probase is a probabilistic KB of millions concepts and their relationships (basically is-a). It was created by mining billions of Web pages and search logs of Microsoft’s Bing⁵ repository using syntactic patterns. The concept KB was then leveraged for text conceptualization to support various text understanding tasks such as clustering of Twitter messages and News titles (Song et al., 2011, 2015), (Song et al., 2015), search query understanding (Wang et al., 2015), short text segmentation (Wang et al., 2014; Hua et al., 2015), and term similarity (Li et al., 2013; Kim et al., 2013).

Despite its effectiveness, the dependency of Probase on syntactic patterns can be a limitation especially for languages other than English. In addition, we expect augmenting and maintaining these syntactic patterns to be costly and labor intensive. We argue that concept embeddings allow simpler and more efficient representations, simply because similarity scoring between concept vectors can be performed using vector arithmetic. While the Probase hierarchy allows only symbolic matching which still suffers data sparsity. On another hand, we spotted some cases where Probase probabilities were atypical⁶. This is due to learning concept categories from a limited set of syntactic patterns which does not cover all concept mention patterns. Concept embeddings relax this problem by leveraging all concept mentions in order to learn the embedding vector and therefore might be utilized to curate such atypical Probase assertions.

Concept/Entity Embeddings: Neural embedding models have been proposed to learn distributed representations of concepts/entities⁷. Song and Roth (2015) proposed using the popular Word2Vec model (Mikolov et al., 2013a) to obtain the embeddings of each concept by averaging the embeddings of the concept’s individual words. For example, the embeddings of *Microsoft Office* would be obtained by averaging the vector of *Microsoft* and the vector of *Office* obtained from the Word2Vec model. Clearly, this method disregards the fact that the semantics of multiword concepts whose composite meaning is different from the semantics of their individual words.

More robust concept and entity embeddings can be learned from the general knowledge about the concept in encyclopedic KB (e.g., its article) and/or from the structure of a hyperlinked KB (e.g., its link graph). Such concept embedding models were proposed by Hu et al. (2015), Li et al. (2016), and Yamada et al. (2016) who all utilize the skip-gram

4. <https://concept.research.microsoft.com>

5. <https://www.bing.com/>

6. $p(\text{Arabic coffee} \mid \text{beverage}) = 0$

7. In this paper we use the terms "concept" and "entity" interchangeably.

learning technique (Mikolov et al., 2013b), but differ in how they define the context of the target concept.

Li et al. (2016) extended the embedding model proposed by Hu et al. (2015) by jointly learning entity and category embeddings from contexts defined by all other entities in the target entity article as well as its category hierarchy in Wikipedia. This method has the advantage of learning embeddings of both entities and categories jointly. However, defining the entity contexts as pairs of the target entity and all other entities appearing in its corresponding article might introduce noisy contexts, especially for long articles. For example, the Wikipedia article for "*United States*" contains links to "*Kindergarten*", "*First grade*", and "*Secondary school*" under the "*Education*" section.

Yamada et al. (2016) proposed a method based on the skip-gram model to jointly learn embeddings of words and entities using contexts generated from surrounding words of the target entity or word. The authors also proposed incorporating *Wikipedia* link graph by generating contexts from all entities with outgoing link to the target entity to better model entity-entity relatedness.

Our models also learn word and concept embeddings jointly. Mapping both words and concepts into the same semantic space allows us to easily measure word-word, word-concept, and concept-concept semantic similarities. In addition, our CRX model (described in Section 3) extends the context of each word/concept by including nearby concept mentions and not only nearby words. Therefore, we better model local contextual information of concepts and words in *Wikipedia*, treated as a textual KB. During training, we generate word-word, word-concept, concept-word, and concept-concept contexts (cf. equation 5). In Yamada et al. (2016) model, concept-concept contexts are generated from *Wikipedia* link graph not from their raw mentions in *Wikipedia* text. In the CCX model, we define concept contexts by all surrounding concepts within a window of fixed size.

Generating contexts from raw text mentions makes our models scalable and *not restricted to hyperlinked encyclopedic textual corpora* only. This facilitates exploiting other free-text corpora with annotated concept mentions (e.g., news stories, scientific publications, medical guidelines...etc). Moreover, our proposed models are *computationally less costly* than Hu et al. (2015) and Yamada et al. (2016) models as they require few hours rather than days to train on similar computing resources.

BoC Densification: Densification of the Bag-of-Concepts (BoC) is the process of converting the sparse BoC into a continuous BoC (CBoC) (aka dense BoC) in order to overcome the BoC sparsity problem. The process requires first mapping each concept into a continuous vector using representation learning. Song and Roth (2015) proposed three different mechanisms for aligning the concepts at different indices given a sparse BoC pair (\mathbf{u}, \mathbf{v}) in order to increase their similarity score.

The *many-to-many* mechanism works by averaging all pairwise similarities. The *many-to-one* mechanism works by aligning each concept in \mathbf{u} with the most similar concept in \mathbf{v} (i.e., its best match). Clearly, the complexity of these two mechanisms is *quadratic*. The third mechanism is the *one-to-one*. It utilizes the Hungarian method in order to find an optimal alignment on a one-to-one basis (Papadimitriou & Steiglitz, 1982). This mechanism performed the best on the task of dataless document classification and was also utilized by Li et al. (2016).

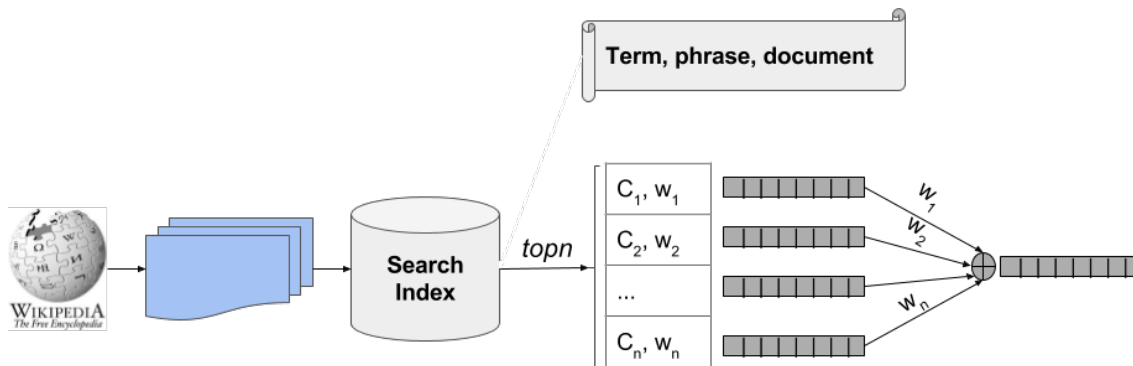


Figure 1: Densification of the bag of concepts vector using weighted average of the learned concept embeddings. The concept space is defined by all Wikipedia article titles. The concept vector is created from the top n hits of searching a *Wikipedia* inverted index with the given text. The weights are the TF-IDF scores from searching *Wikipedia*.

However, the Hungarian method is a combinatorial optimization algorithm whose complexity is *polynomial*. Our proposed densification mechanism is more efficient than these three mechanisms as its complexity is *linear* with respect to the number of *nonzero* elements in the BoC. Additionally, it is simpler as it does not require tuning a cutoff threshold for the minimum similarity between two aligned concepts as in previous work. Figure 1 shows a schematic diagram of our efficient densification mechanism applied to a BoC generated from a *Wikipedia* inverted index. We simply perform weighted average of the individual concept vectors in the obtained BoC where concept weights correspond to the TF-IDF scores from searching *Wikipedia*.

3. Learning Concept Embeddings

A main objective of learning concept embeddings is to overcome the inherent problem of *data sparsity* associated with the BoC representation. Here we try to learn continuous concept vectors by building upon the skip-gram embedding model (Mikolov et al., 2013b). In the conventional skip-gram model, a set of contexts are generated by sliding a context window of predefined size over sentences of a given text corpus. Vector representation of a target word is learned with the objective to maximize the ability of predicting surrounding words of that target word.

Formally, given a training corpus of V words $\omega_1, \omega_2, \dots, \omega_V$. The skip-gram model aims to maximize the average log probability:

$$\frac{1}{V} \sum_{i=1}^V \sum_{-s \leq j \leq s, j \neq 0} \log p(\omega_{i+j} | \omega_i) \quad (3)$$

where s is the context window size, ω_i is the target word, and ω_{i+j} is a surrounding context word. The softmax function is used to estimate the probability $p(\omega_O|\omega_I)$ as follows:

$$p(\omega_O|\omega_I) = \frac{\exp(\mathbf{v}_{\omega_O}^\top \mathbf{u}_{\omega_I})}{\sum_{\omega=1}^W \exp(\mathbf{v}_w^\top \mathbf{u}_{\omega_I})} \quad (4)$$

where \mathbf{u}_{ω_I} and \mathbf{v}_{ω_O} are the input and output vectors respectively and W is the vocabulary size. Mikolov et al. (2013b) proposed hierarchical softmax and negative sampling as efficient alternatives to approximate the softmax function which becomes computationally intractable when W becomes huge.

Our approach genuinely learns distributed concept representations by generating concept contexts from *mentions* of those concepts in large encyclopedic KBs such as *Wikipedia*. Utilizing such annotated KBs eliminates the need to manually annotate concept mentions and thus comes at no cost.

3.1 Concept Raw Context Model (CRX)

In this model, we jointly learn the embeddings of both words and concepts. First, all concept mentions are identified in the given corpus. Second, contexts are generated for both words and concepts from other surrounding words and other surrounding concepts as well. After generating all the contexts, we use the skip-gram model to jointly learn the embeddings of words and concepts. Formally, given a training corpus of V words $\omega_1, \omega_2, \dots, \omega_V$, we iterate over the corpus identifying words and concept mentions and thus generating a sequence of T tokens t_1, t_2, \dots, t_T where $T < V$ (as multiword concepts will be counted as one token). Afterwards we train the a skip-gram model aiming to maximize:

$$\frac{1}{T} \sum_{i=1}^T \sum_{-s \leq j \leq s, j \neq 0} \log p(t_{i+j}|t_i) \quad (5)$$

where as in the conventional skip-gram model, s is the context window size. In this model, t_i is the target token which would be either a word or a concept mention, and t_{i+j} is a surrounding context word or concept mention.

This model is different from Yamada et al. (2016)’s anchor context model in three aspects: 1) while generating target concept contexts, we utilize not only surrounding words but also other surrounding concepts, 2) our model aims to maximize $p(t_{i+j}|t_i)$ where t could be a word or a concept, while Yamada et al. (2016) model maximizes $p(\omega_{i+j}|e_i)$ where e_i is the target concept/entity (see Yamada et al. (2016) Eq. 6), and 3) in case t_i is a concept, our model captures all the contexts in which it appeared, while Yamada et al. (2016) model generates for each entity one context of s previous and s next words. We hypothesize that considering both concepts and individual words in the optimization function generates more robust embeddings.

3.2 Concept-Concept Context Model (CCX)

Inspired by the distributional hypothesis (Harris, 1954), in this model, we hypothesize that: *”similar concepts tend to appear in similar conceptual contexts”*. In order to test this hypothesis, we propose learning concept embeddings by training a skip-gram model

Sentence	CRX Contexts	CCX Contexts
<i>Larry Page</i> is the co-founder of <i>Google</i> which is headquartered in <i>Menlo Park CA</i>	< <i>Larry Page</i> , co-founder> <co-founder, <i>Google</i> > < <i>Google</i> , headquartered> <headquartered, <i>Menlo Park CA</i> >	< <i>Larry Page</i> , <i>Google</i> > < <i>Larry Page</i> , <i>Menlo Park CA</i> > < <i>Google</i> , <i>Menlo Park CA</i> >
<i>Bill Gates</i> is the co-founder of <i>Microsoft</i> which is headquartered in <i>Redmond WA</i>	< <i>Bill Gates</i> , co-founder> <co-founder, <i>Microsoft</i> > < <i>Microsoft</i> , headquartered> <headquartered, <i>Redmond WA</i> >	< <i>Bill Gates</i> , <i>Microsoft</i> > < <i>Bill Gates</i> , <i>Redmond WA</i> > < <i>Microsoft</i> , <i>Redmond WA</i> >
<i>Google</i> is headquartered in <i>Menlo Park CA</i> and was co-founded by <i>Larry Page</i>	< <i>Google</i> , headquartered> <headquartered, <i>Menlo Park CA</i> > < <i>Menlo Park CA</i> , co-founded> <co-founded, <i>Larry Page</i> >	< <i>Google</i> , <i>Menlo Park CA</i> > < <i>Google</i> , <i>Larry Page</i> > < <i>Menlo Park CA</i> , <i>Larry Page</i> >

Table 2: Example three sentences along with sample contexts generated from CRX and CCX. Contexts are generated with a context window of length 3.

on contexts generated solely from concept mentions. As in the CRX model, we start by identifying all concept mentions in the given corpus. Then, contexts of target concept are generated from surrounding concepts only. Formally, given a training corpus of V words $\omega_1, \omega_2, \dots, \omega_V$. We iterate over the corpus identifying concept mentions and thus generating a sequence of C concept tokens c_1, c_2, \dots, c_C where $C < V$. Afterwards we train the skip-gram model aiming to maximize:

$$\frac{1}{C} \sum_{i=1}^C \sum_{-s \leq j \leq s, j \neq 0} \log p(c_{i+j} | c_i) \quad (6)$$

where s is the context window size, c_i is the target concept, and c_{i+j} is a surrounding concept mention within s mentions.

This model is different from Li et al. (2016) and Hu et al. (2015) as they define the context of a target concept by all the other concepts which have an outgoing link from the concept’s corresponding article in *Wikipedia*.

Formally, given an article about concept c_t containing other concepts (c_1, c_2, \dots, c_n) , Li et al. (2016), Hu et al. (2015) create context pairs in the form (c_t, c_i) , $1 \leq i \leq n$. Thus context size is limited to only 2 concepts.

Clearly, some of these concepts might be irrelevant especially for very long articles which cite hundreds of other concepts. Our CCX model, alternatively, learns concept semantics from surrounding concepts and not only from those that are cited in its article. We also extend the context window beyond pairs of concepts (based on s in equation 6) allowing more influence to other nearby concepts.

3.3 CRX vs. CCX

One of the advantages of the CCX model over the CRX model is its computational efficiency during learning. On the other hand, the CCX model vocabulary is *limited to the corpus concepts* (all *Wikipedia* articles in our case), while the CRX model vocabulary is defined by all *unique concepts+words* in *Wikipedia*.

Another distinct property of the CCX model is its emphasis on concept-concept *relatedness rather than similarity* (as we will detail more in the experiments section). The CCX model by looking only at surrounding concept mentions while learning, is able to generate contexts containing more diverse but related concepts. On the other hand, the CRX model which jointly learns the embeddings of words and concepts puts more *emphasis on similarity* by leveraging the full contextual information of words and concepts while learning.

To better illustrate this difference, consider a sample of the contexts generated from CRX and CCX in Table 2 using a sliding window of length 3. As we can notice, the CRX contexts of "Google" and "Microsoft" are somewhat similar containing words like "headquartered" and "co-founder". This causes the model to learn similar vectors for these two concepts. On the other hand, the CCX contexts of "Google" and "Microsoft" do not share any similarities⁸, rather we can see that "Google" has similar contexts to "Larry Page" as both has "Menlo Park CA" in their contexts, causing the model to learn similar embeddings for these two related concepts.

3.4 Training

We utilize a recent *Wikipedia* dump of August 2016⁹, which has about 7 million articles. We extract articles plain text discarding images and tables. We also discard "References" and "External links" sections (if any). We pruned both articles not under the main namespace and pruned all redirect pages as well. Eventually, our corpus contained about 5 million articles in total.

We preprocess each article replacing all its references to other *Wikipedia* articles with their corresponding page IDs. In case any of the references is a title of a redirect page, we use the page ID of the original page to ensure that all concept mentions are normalized.

Following Mikolov et al. (2013b), we utilize negative sampling to approximate the softmax function by replacing every $\log p(\omega_O|\omega_I)$ term in the softmax function (equation 4) with:

$$\log \sigma(\mathbf{v}_{\omega_O}^T \mathbf{u}_{\omega_I}) + \sum_{s=1}^k \mathbb{E}_{\omega_s \sim P_n(w)} [\log \sigma(-\mathbf{v}_{\omega_s}^T \mathbf{u}_{\omega_I})] \quad (7)$$

where k is the number of negative samples drawn for each word and $\sigma(x)$ is the sigmoid function ($\frac{1}{1+e^{-x}}$). In the case of the CRX model ω_I and ω_O would be replaced with t_i and t_{i+j} respectively. And in the case of the CCX model ω_I and ω_O would be replaced with c_i and c_{i+j} respectively.

For both the CRX & CCX models we use a context window of size 9 and a vector of 500 dimensions. We train the skip-gram model for 10 iterations using 12-core machine with 64GB of RAM. The CRX model took ~ 15 hours to train for a total of ~ 12.7 million tokens. The CCX model took ~ 1.5 hours to train for a total of ~ 4.5 million concepts.

3.5 BoC Densification

As we mentioned in the related work section, the current mechanisms for BoC densification are inefficient as their complexity is at least quadratic with respect to the number of nonzero

8. This is an illustrative example and doesn't imply the two concepts will have totally dissimilar vectors.

9. <http://dumps.wikimedia.org/enwiki/>

elements in the BoC vector. Here, we propose simple and efficient vector aggregation method to obtain fully continuous BoC vectors (CBoC) in linear time. Our mechanism works by performing a weighted average of the individual concept vectors in a given BoC. This operation has two advantages. First, it *scales linearly* with the number of nonzero dimensions in the BoC vector. Secondly, it produces a fully dense BoC vector of *fixed size* representing the semantics of the original concepts and *considering their weights*. Formally, given a sparse BoC vector $\mathbf{s} = \{(c_1, w_1), \dots, (c_{|s|}, w_{|s|})\}$, where w_i is weight of concept c_i ¹⁰. We can obtain the dense representation of \mathbf{s} as in equation 8:

$$\mathbf{s}_{dense} = \frac{\sum_{i=1}^{|\mathbf{s}|} w_i \cdot \mathbf{u}_{c_i}}{\sum_{i=1}^{|\mathbf{s}|} w_i} \quad (8)$$

where \mathbf{u}_{c_i} is the vector of concept c_i . Once we have this dense BoC vector, we can apply the cosine measure to compute the similarity between a pair of dense BoC vectors.

As we can notice, this weighted average is done *once* for any given BoC vector. Other mechanisms that rely on concept alignment (Song & Roth, 2015), require *realignment* every time a pair of BoC vectors are compared. Our approach improves the *efficiency* especially in the context of dataless document classification with large number of classes. Using our densification mechanism, we apply the weighted average for the BoC of each category and for each instance document once.

Interestingly, our densification mechanism allows us to densify the sparse BoC vector using only the *top few dimensions*. As we will show in the experiments, we can get *near-best* results using these few dimensions compared to densifying with all the dimensions in the original sparse vector. This property reduces the cost of obtaining a BoC vector with a few hundred dimensions in the first place.

4. Text Conceptualization Applications

Concept-based representations have many applications in computational linguistics, information retrieval, and knowledge modeling. Such representations are able to capture the semantics of a given text by either identifying concept mentions in that text, transforming the text into a concept space, or both (Wang & Wang, 2016). Thereafter, many cognitive tasks that require huge background and real-world knowledge are facilitated by leveraging the conceptual representations. We describe some of these tasks in this section, and provide empirical evaluation of our our concept embedding models on such tasks in the next section.

4.1 Concept/Entity Relatedness

Entity relatedness has been recently used to model *entity coherence* in many named entity linking and disambiguation systems (Witten & Milne, 2008; Milne & Witten, 2008; Hoffart et al., 2012; Ceccarelli et al., 2013; Huang et al., 2015; Hu et al., 2015; Yamada et al., 2016). In entity search, Hu et al. (2015) utilized entity relatedness score to *rank* candidate entities based on their relatedness to the search query entities. Also, entity embeddings have proved more efficient and effective for measuring entity relatedness over traditional relatedness measures which use link analysis. Formally, given a entity pair (e_i, e_j) , their

10. The weights are the TF-IDF scores from searching *Wikipedia*.

relatedness score is evaluated as $rel(e_i, e_j) = Sim(\mathbf{u}_{e_i}, \mathbf{u}_{e_j})$, where Sim is a similarity function (e.g., *cosine*), and \mathbf{u}_e is the embeddings of entity e .

4.2 Concept Learning

Concept learning is a cognitive process which involves classifying a given concept/entity to one or more candidate categories (e.g., *milk* as *beverage*, *dairy product*, *liquid...etc*). This process is also known as *concept categorization*¹¹ (Li et al., 2016). Automated concept learning gains its importance in many knowledge modeling tasks such as knowledge base *construction* (discovering new concepts), *completion* (inferring new relationships between concepts), and *curation* (removing noisy or assessing weak relationships). Similar to Li et al. (2016), we assign a given concept to a target category using Rocchio classification (Rocchio, 1971), where the centroid of each category is set to the category’s corresponding embedding vector. Formally, given a set of n candidate concept categories $G = \{g_1, \dots, g_n\}$, a sample concept c , an embedding function f , and a similarity function Sim , then c_i is assigned to category g_* such that $g_* = arg \max_i Sim(f(g_i), f(c))$. Here, the embedding function f would always map the given concept to its vector.

4.3 Dataless Classification

Chang et al. (2008) proposed dataless document classification as a learning *protocol* to perform text categorization without the need for labeled data to train a classifier. Given only label names and few descriptive keywords of each label, classification is performed *on the fly* by mapping each label into a BoC representation using ESA (Gabrilovich & Markovitch, 2007). Likewise, each data instance is mapped into the same BoC semantic space and assigned to the most similar label using a proper similarity measure such as *cosine*. Formally, given a set of n labels $L = \{l_1, \dots, l_n\}$, a text document d , a BoC mapping model f , and a similarity function Sim . First we conceptualize the each l_i and the document d by applying f on them, which will produce sparse BoC vectors s_{l_i} and s_d respectively. Then we densify the vectors as in equation 8 producing $s_{dense_{l_i}}$ and s_{dense_d} respectively. Finally d is assigned to label l_* such that $l_* = arg \max_i Sim(s_{dense_{l_i}}, s_{dense_d})$.

In the context of dataless classification, Chang et al. (2008) and Song and Roth (2014) used bootstrapping in order to improve the classification performance without the need for labeled data. The basic idea is to start from target labels as the initial training samples, train a classifier, and *iteratively* add to the training data those samples which the classifier is *most confident* until no more samples to be classified. The results of dataless classification with bootstrapping were competitive to supervised classification with many training examples.

In this paper, we extend the use of bootstrapping to the concept learning task as well. In concept learning we start with the vectors of target category concepts as a prototype view upon which categorization decisions are made (e.g., $vec(bird)$, $vec(mammal)$...etc). We leverage bootstrapping by *iteratively updating* this prototype view with the vectors of concept instances we are most confident. For example, if *”deer”* is closest to *”mammal”* than any other instance in the dataset, then we update the definition of *”mammal”* by performing $vec(mammal) += vec(deer)$, and repeat the same operation for other categories as well.

11. In this paper, we use concept learning and concept categorization interchangeably

Algorithm 1: Classification + Bootstrapping

Input: $\mathbf{U} = \{(l_1, \mathbf{u}_{l_1}), \dots, (l_n, \mathbf{u}_{l_n})\}$: labels + embeddings
 $\mathbf{D} = \{(d_1, \mathbf{v}_{d_1}), \dots, (d_m, \mathbf{v}_{d_m})\}$: instances + embeddings
 N : number of bootstrap instances
Result: $\mathbf{L} = \{\dots, (d_i, l_j), \dots\}$: label assignment for each instance

```
1 repeat
2   candidates  $\leftarrow \{l_1 : \phi, \dots, l_n : \phi\}$ 
3   foreach  $(d, \mathbf{v}_d) \in \mathbf{D}$  do
4      $d_{max.sim} = 0$ 
5      $d_{max.label} = null$ 
6     foreach  $(l, \mathbf{u}_l) \in \mathbf{U}$  do
7        $sim_l = Sim(\mathbf{v}_d, \mathbf{u}_l)$ 
8       if  $sim_l > d_{max.sim}$  then
9          $d_{max.sim} = sim_l$ 
10         $d_{max.label} = l$ 
11      end
12    end
13    add  $(d, d_{max.sim})$  to candidates[l]
14  end
15  foreach  $(l, candidates_l) \in candidates.items$  do
16    repeat
17       $score_{max} = 0$ 
18       $d_{max} = null$ 
19      foreach  $(d, score_d) \in candidates_l$  do
20        if  $score_d > score_{max}$  then
21           $score_{max} = score_d$ 
22           $d_{max} = d$  ▷ most similar instance so far
23        end
24      end
25      add  $(d_{max}, l)$  to  $\mathbf{L}$  ▷ assign class label
26       $\mathbf{u}_l \leftarrow \mathbf{u}_l + \mathbf{v}_d$  ▷ bootstrap label embedding
27      remove  $d$  from candidates_l
28      remove  $d$  from  $\mathbf{D}$ 
29    until  $N$  highest scored instances added
30  end
31 until  $\mathbf{D} = \phi$  ▷ no more instances to classify
```

This way, we *adapt* the initial prototype view to better match the specifics of the given data. Although bootstrapping is a time consuming process, we argue that, using dense vectors for representing concepts makes bootstrapping more appealing. As updating the category vector with an instance vector could be performed through optimized vector arithmetic which is available in most modern machines. Algorithm 1 presents the pseudocode for performing dataless classification and concept categorization with bootstrapping. In our implementation, we bootstrap the category vector with vectors of the most similar \mathbf{N} instances at a time. Another implementation option might be defining a threshold and bootstrapping using vectors of \mathbf{N} instances if their similarity score exceed that threshold. In the experiments, we set $\mathbf{N}=1$.

5. Experiments

5.1 Entity Semantic Relatedness

We evaluate the "goodness" of our concept embeddings on measuring entity semantic relatedness as an intrinsic evaluation.

Method	IT Companies	Celebrities	TV Series	Video Games	Chuck Norris	All
WLM	0.721	0.667	0.628	0.431	0.571	0.610
CombIC	0.644	0.690	0.643	0.532	0.558	0.624
ExRel	0.727	0.643	0.633	0.519	0.477	0.630
KORE	0.759	0.715	0.599	0.760	0.498	0.698
CRX	0.644	0.592	0.511	0.641	0.495	0.586
CCX	0.788	<u>0.694</u>	0.696	<u>0.708</u>	0.573	0.714

Table 3: Evaluation of concept embeddings for measuring entity semantic relatedness using Spearman rank-order correlation (ρ). Overall, the CCX model gives the best results outperforming all other models. It comes 1st on 3 categories (bold), and 2nd on the other two (underlined).

5.1.1 DATASET

We use the KORE dataset created by Hoffart et al. (2012). It consists of 21 main entities from four domains: IT companies, Hollywood celebrities, video games, and television series. For each of these entities, 20 other candidate entities were selected and manually ranked based on their relatedness score based on human judgements. As in previous studies, we report the Spearman rank-order correlation (ρ) (Zwillinger & Kokoska, 1999) which assesses how the automated ranking of candidate entities based on their relatedness score matches the ranking we obtain from human judgements.

5.1.2 COMPARED SYSTEMS

We compare our models with four previous methods:

1. **KORE** (Hoffart et al., 2012) which measure entity relatedness by firstly representing entities as sets of weighted keyphrases and then computing relatedness using different measures such as keyphrase vector cosine similarity and keyphrase overlap relatedness.
2. **WLM** introduced by Witten and Milne (2008) who proposed a Wikipedia Link-based Measure (WLM) as a simple mechanism for modeling the semantic relatedness between *Wikipedia* entities. The authors utilized *Wikipedia* link structure under the assumption that related entities would have similar incoming links.
3. **Exclusivity-based Relatedness (ExRel)** introduced by Hulpus et al. (2015) who proposed this measure under the assumption that not all instances of a given relation type should be equally weighted. Specifically, the authors hypothesized that the relatedness score between two concepts should be higher if each of them is related through the same relation type to fewer other concepts in the employed KB link graph.
4. **Combined Information Content (CombIC)** introduced by Schuhmacher and Ponzetto (2014) who compute the relatedness score using a graph edit distance measure on the *DBpedia KB*.

Entity	CRX	CCX	Ground Truth (GT)
Google	Yahoo! (9) Apple Inc. (12) Bing (search engine) (7)	<u>Larry Page</u> (1) <u>Sergey Brin</u> (2) Yahoo! (9)	Larry Page Sergey Brin Google Maps
Leonardo DiCaprio	Kate Winslet (4) Steven Spielberg (9) Tobey Maguire (7)	Tobey Maguire (7) Kate Winslet (4) <u>Titanic (1997 film)</u> (2)	Inception (film) Titanic (1997 film) Frank Abagnale
Mad Men	The Sopranos (15) <u>Matthew Weiner</u> (1) <u>Jon Hamm</u> (2)	<u>Matthew Weiner</u> (1) <u>Jon Hamm</u> (2) Todd London (4)	Matthew Weiner Jon Hamm Alan Taylor (director)
Guitar Hero (video game)	Frequency (video game) (10) Rock Band (video game) (6) Harmonix Music Systems (1)	<u>Harmonix Music Systems</u> (1) <u>WaveGroup Sound</u> (3) <u>RedOctane</u> (1)	Harmonix Music Systems RedOctane WaveGroup Sound

Table 4: Top-3 rated entities from CRX & CCX models on sample entities from the 4 domains compared to the ground truth. We can notice high agreement between CCX model ranks and the ground truth ranks (in brackets). The CRX model top rated entities has lower ranks than ground truth ranks causing relatively low correlation scores.

5.1.3 RESULTS

Table 3 shows the Spearman (ρ) correlation scores of the CRX and CCX model compared to previous models. As we can notice the CCX model achieves the best overall performance on the five domains combined exceeding its successor KORE by 1.6%. The CRX model on the other hand came last on this task.

In order to better understand these results, we looked at rankings of individual entities from each domain to see how they compare to the ground truth. Table 4 shows the top-3 rated entities from each model on sample entities from the four domains. As we can notice, the ground truth assigns high rank to related rather than similar entities. For example, relatedness of "Google" to "Larry Page" is ranked 1st, while to "Yahoo!" is ranked 9th, and to "Apple Inc." is ranked 12th. As the CCX model emphasizes semantic relatedness over similarity, it has high overlap in the top-3 entities with the ground truth (underlined entities). On the other hand, the CRX model predictions are actually meaningful when it comes to functional and topical similarity. As we can notice, it assigns high ranks of "Google" to other companies ("Yahoo!", "Apple Inc."), of "Leonardo DiCaprio" to other celebrities ("Tobey Maguire"), and "Mad Men" to other TV series ("The Sopranos"), and of "Guitar Hero" to other video games ("Frequency", "Rock Band"). However, all these highly ranked entities by CRX have relatively low rankings in the ground truth (given in brackets). This caused the correlation score to be much lower than what we obtained from the CCX model.

The results indicate that, the CCX model could be more appropriate in applications where relatedness and topical diversity are more desired than topical and functional coherence where the CRX model would be more appealing.

Dataset/Instances Method	Battig (83)	DOTA-single (300)	DOTA-mult (150)	DOTA-all (450)
WE _{Senna}	0.44	0.52	0.32	0.45
WE _{Mikolov}	0.74	0.72	0.67	0.72
TransE ₁	0.66	0.72	0.69	0.71
TransE ₂	0.75	0.80	0.77	0.79
TransE ₃	0.46	0.55	0.52	0.54
CE	0.79	0.89	0.85	0.88
HCE	0.87	0.93	0.91	0.92
CCX	0.72	0.90	0.80	0.87
+bootstrap	0.81	0.91	0.85	0.87
CRX	0.83	0.91	0.88	0.90
+bootstrap	0.89	0.98	0.95	0.97

Table 5: Accuracy of concept categorization. The CRX model with bootstrapping gives the best results outperforming all other models.

5.2 Concept Categorization

This task can be viewed as both intrinsic and extrinsic. It is intrinsic because a *good* embedding model would generate clusters of concepts belonging to the same category, and optimally place the category vector at the center of its instances vectors. On another hand, it is extrinsic as the embedding model could be used to generate a concept KB of is-a relationships with confidence scores, similar to *Probase* (Wu et al., 2012). The model could even be used to curate and/or assert the facts in *Probase*.

5.2.1 DATASETS

As in Li et al. (2016), we utilize two benchmark datasets: 1) Battig test (Baroni & Lenci, 2010), which contains 83 single word concepts (e.g., *cat, tuna, spoon..etc*) belonging to 10 categories (e.g., *mammal, fish, kitchenware..etc*), and 2) DOTA which was created by Li et al. (2016) from *Wikipedia* article titles (entities) and category names (categories). DOTA contains 300 single-word concepts (DOTA-single) (e.g., *coffee, football, semantics..etc*), and (150) multiword concepts (DOTA-mult) (e.g., *masala chai, table tennis, noun phrase...etc*). Both belong to 15 categories (e.g., *beverage, sport, linguistics..etc*). Performance is measured in terms of the ability of the system to assign concept instances to their correct categories.

5.2.2 COMPARED SYSTEMS

We compare our models to various word, entity, and category embedding methods as described in Li et al. (2016) including:

1. **Word embeddings:** Collobert et al. (2011) model (WE_{Senna}) trained on *Wikipedia*. Here vectors of multiword concepts are obtained by averaging their individual word vectors.

2. ***MWEs embeddings***: Mikolov et al. (2013b) model ($WE_{Mikolov}$) trained on *Wikipedia*. This model jointly learns single and multiword embeddings where MWEs are identified using corpus statistics.
3. ***Entity-category embeddings***: which include Bordes et al. (2013) embedding model (TransE). This model utilizes relational data between entities in a KB as triplets in the form (entity,relation,entity) to generate representations of both entities and relationships. Li et al. (2016) implemented three variants of this model (TransE₁, TransE₂, TransE₃) to generate representations for entities and categories jointly. Two other models introduced by Li et al. (2016) are CE and HCE. CE generates embeddings for concepts and categories using category information of *Wikipedia* articles. HCE extends CE by incorporating *Wikipedia*'s category hierarchy while training the model to generate concept and category vectors.

5.2.3 RESULTS

We report the accuracy scores of concept categorization¹² in Table 5. Accuracy is calculated by dividing the number of correctly classified concepts by the total number of concepts in the given dataset. Scores of all other methods are obtained from Li et al. (2016). As we can see in Table 5, the CRX model comes second after the HCE on all datasets. While the CCX model performance is much less than CRX. With bootstrapping, the CCX model performance improves on both datasets. CRX with bootstrapping outperforms all other models by significant percentages. These results show that learning concept embeddings from concept mentions is actually different from training the skip-gram model on phrases or multiword expressions. This is clear from the significant performance gains we get from the CRX and CCX models compared to $WE_{Mikolov}$ which was trained using skip-gram on phrases. Additionally, the results demonstrate the efficacy of our models which simply learn concept embeddings from concept mentions in free-text corpus compared to the more complex models which require category or relational information such as TransE, CE, and HCE.

5.3 Dataless Classification

In this experiment, we evaluate the effectiveness of our concept embedding models on the dataless document classification task as an extrinsic evaluation. We demonstrate through empirical results the efficiency and effectiveness of our proposed BoC densification scheme which helps obtaining better classification results compared to the original sparse BoC representation.

5.3.1 DATASET

We use the 20-newsgroups dataset (20NG) (Lang, 1995) which is commonly used for benchmarking text classification algorithms. The dataset contains 20 categories each has ~ 1000 news posts. We obtained the BoC representations using ESA from Song and Roth (2014) who utilized a Wikipedia index containing pages with 100+ words and 5+ outgoing links

12. From a multi-class classification perspective, the accuracy scores would be equivalent to the clustering purity score as reported in Li et al. (2016).

Top-level	Low-level
Sport	Hockey, Baseball, Autos, Motorcycles
Politics	Guns, Mideast, Misc
Religion	Christian, Atheism, Misc

Table 6: 20NG category mappings

Method	Hockey x Baseball		Autos x Motorcycles		Guns x Mideast x Misc	
ESA	94.60	@425	72.70	@325	70.00	@500
CCX (equal)	94.60	@20	-	-	70.33	@60
CRX (equal)	94.60	@60	73.10	@4	70.00	@7
WE_{max}	86.85	@65	76.15	@375	72.20	@300
WE_{hung}	95.20	@325	73.75	@300	71.70	@275
CCX (best)	95.10	@125	69.70	@7	72.47	@250
+bootstrap	95.90	@450	74.25	@12	77.43	@5
CRX (best)	95.65	@425	79.20	@14	73.40	@70
+bootstrap	95.90	@350	73.25	@12	77.03	@10

Table 7: Evaluation results of dataless document classification of fine-grained classes measured in micro-averaged F1 along with # of dimensions (concepts) in the BoC at which corresponding performance is achieved.

to create ESA mappings of 500 dimensions for both the categories and news posts of the 20NG. We designed two types of classification tasks: 1) fine-grained classification involving closely related classes such as *Hockey* vs. *Baseball*, *Autos* vs. *Motorcycles*, and *Guns* vs. *Mideast* vs. *Misc*, and 2) coarse-grained classification involving top-level categories such as *Sport* vs. *Politics* and *Sport* vs. *Religion*. The top-level categories are created by combining instances of the fine-grained categories which are shown in Table 6.

5.3.2 COMPARED SYSTEMS

We compare our models to three previous methods:

1. **ESA** which computes the cosine similarity between target labels and instance documents using the sparse BoC vectors.
2. **WE_{max}** & **WE_{hung}** which were proposed by Song and Roth (2015) for BoC den-sification using embeddings obtained from Word2Vec. As the authors reported, we fix the minimum similarity threshold to 0.85. WE_{max} finds the best match for each concept, while WE_{hung} utilizes the Hungarian algorithm to find the best concept-concept alignment on one-to-one basis. Both mechanisms have polynomial degree time complexity.

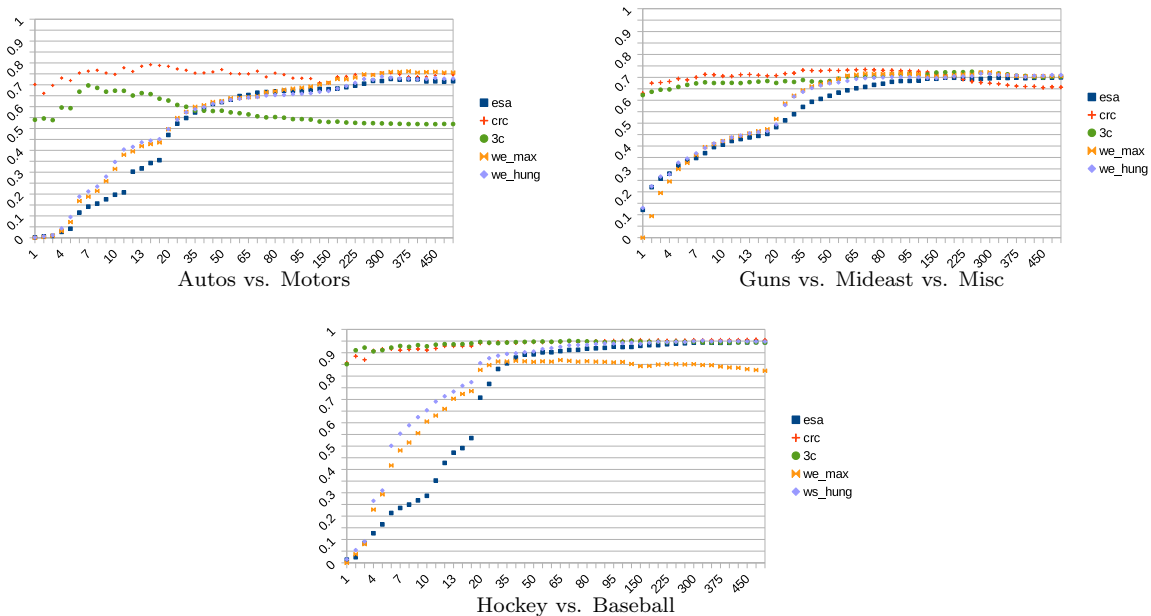


Figure 2: micro-averaged F1 scores of fine-grained classes when varying the # of BoC dimensions.

5.3.3 RESULTS

Table 7 presents the results of fine-grained dataless classification measured in micro-averaged F1. As we can notice, ESA achieves its peak performance with a few hundred dimensions of the sparse BoC vector. Using our densification mechanism, both the CRX & CCX models achieve equal performance to ESA at many fewer dimensions. Densification using the CRX model embeddings gives the best F1 scores on the three tasks. Interestingly, the CRX model improves the F1 score by $\sim 7\%$ using only 14 concepts on *Autos vs. Motorcycles*, and by $\sim 3\%$ using 70 concepts on *Guns vs. Mideast vs. Misc*. The CCX model, still performs better than ESA on 2 out of the 3 tasks. Both WE_{max} and WE_{hung} improve the performance over ESA but not as our CRX model.

When we applied bootstrapping, the performance of the CCX model improved slightly on *Hockey vs. Baseball*, but significantly ($\sim 5\%$) on the other two tasks achieving best performance on the third task with just 5 concepts. Bootstrapping with the CRX model has a similar effect to the CCX model except for *Autos vs. Motorcycles* where performance degraded significantly. To better understand this behavior, we analyzed the results as bootstrapping progresses at 14 concepts like CRX (best). We noticed that, at the very early iterations of Algorithm 1, many instances belonging to *Autos* were closer to *Motorcycles* with similarity scores between 0.90-0.95. And when using those instances to bootstrap *Motorcycles*, they caused *topic drift* moving *Motorcycles*'s centroid toward *Autos*, and eventually causing relatively lower accuracy scores.

In order to better illustrate the robustness of our densification mechanism when varying the number of BoC dimensions, we measured F1 scores of each task as a function of the

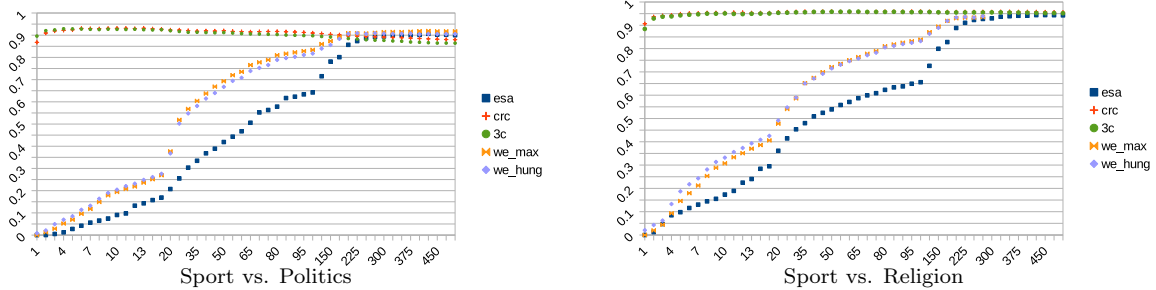


Figure 3: micro-averaged F1 scores of coarse-grained classes when varying the # of concepts (dimensions) in the BoC from 1 to 500.

Method	Sport x Politics		Sport x Religion	
ESA	90.63	@425	94.39	@450
CCX (equal)	92.04	@2	95.11	@6
CRX (equal)	90.99	@2	94.81	@5
WE_{max}	91.89	@425	93.99	@425
WE_{hung}	90.89	@275	94.16	@450
CCX (best)	92.89	@4	95.86	@60
+bootstrap	93.20	@10	95.13	@225
CRX (best)	93.12	@13	95.91	@95
+bootstrap	92.96	@13	95.53	@70

Table 8: Evaluation results of dataless document classification of coarse-grained classes measured in micro-averaged F1 along with # of dimensions (concepts) at which corresponding performance is achieved.

number of BoC dimensions used for densification. As we see in Figure 2, with *one* concept we can achieve high F1 scores compared to ESA which achieves *zero* or very low score. Moreover, *near-peak performance* is achievable with the top 50 or less dimensions. We can also notice that, as we increase the number of dimensions, both WE_{max} and WE_{hung} densification methods have the same undesired monotonic pattern like ESA. Actually, the imposed threshold by these methods does not allow for full dense representation of the BoC vector and therefore at low dimensions we still see low overall F1 score. Our proposed densification mechanism besides its low cost, produces fully densified representations allowing good similarities at low dimensions.

Results of coarse-grained classification are presented in Table 8. Classification at the top level is easier than the fine-grained level. Nevertheless, as with fine-grained classification, ESA still peaks with a few hundred dimensions of the sparse BoC vector. Both the CRX & CCX models achieve equal performance to ESA at very few dimensions (≤ 6). Densification

Method	Hockey x Baseball		Autos x Motorcycles		Guns x Mideast x Misc	
SVM	91.61	@85%	79.25	@20%	77.56	@25%
CCX (best)	95.90	@450	74.25	@12	77.43	@5
CRX (best)	95.90	@350	79.20	@14	77.03	@10

Table 9: Evaluation results of dataless document classification of coarse-grained classes vs. supervised classification with SVM measured in micro-averaged F1 along with # of dimensions (concepts) for CRX & CCX or % of labeled samples for SVM at which corresponding performance is achieved.

using the CRX model embeddings still performs the best on both tasks. Interestingly, the CCX model gives very close F1 scores to the CRX model at less dimensions (@4 with *Sport vs. Politics*, and @60 with *Sport vs. Religion*) indicating its competitive advantage when training computational cost is a decisive criteria. The CCX model, still performs better than ESA, WE_{max} , and WE_{hung} on both tasks.

Bootstrapping did not improve the results on this task significantly (if any). As we can notice in Table 8, the accuracy without bootstrapping is already high indicating that the initial prototype vector (centroid) of each class is representative enough of the instances to be classified.

Figure 3 shows F1 scores of coarse-grained classification when varying the # of BoC dimensions used for densification. The same pattern of achieving *near-peak performance* at very few dimensions recur with the CRX & CCX models. ESA using the sparse BoC vectors achieves low F1 up until few hundred dimensions are considered. Even with the costly WE_{max} and WE_{hung} densifications, performance sometimes decreases.

5.3.4 DATALESS VS. SUPERVISED CLASSIFICATION

We performed a pilot experiment to demonstrate the value of the dataless classification scheme in the absence or difficulty of obtaining labeled data for training a supervised classifier. For this purpose, we used a Support Vector Machine (SVM) classifier with a linear kernel, leveraging the scikit-learn machine learning library (Pedregosa et al., 2011) to perform classification of fine-grained classes (cf. Table 6). We trained the SVM classifier with labeled samples ranging between 10% to 90% of the total number of samples for each task and evaluate the performance on the rest. The results in Table 9 shows the % of labeled sampled needed for training SVM in order to achieve equal performance to dataless classification with CRX and CCX. As we can notice, with *Hockey vs Baseball*, the SVM classifier can't reach the same performance as either models and peaks when trained on 85% (~1700 samples) of the data. With the *Autos vs Motorcycles* and *Guns vs Mideast vs Misc*, the SVM classifier achieves equal performance when trained on 20% (~400 samples) and 25% (~750 samples) of the data respectively. These results demonstrate the competitiveness of our models to supervised classification even when training data is available. And its superiority when training data is scarce.

Concept	Concept Raw Context (CRX)	Concept-Concept Context (CCX)
YouTube	Vevo Facebook SoundCloud Vimeo Viral video	Viral video Vimeo Vevo Video blog Dailymotion
Harvard University	Yale University Princeton University Brown University Columbia University Boston University	Harvard Kennedy School Cambridge, Massachusetts Harvard College Radcliffe College Harvard Society of Fellows
Black hole	Neutron star Accretion disk Primordial black hole Supermassive black hole Event horizon	Event horizon Neutron star Gravitational singularity Wormhole Hawking radiation
X-Men: Days of Future Past	X-Men: Apocalypse X-Men: First Class Deadpool (film) Avengers: Age of Ultron X-Men: The Last Stand	X-Men: Apocalypse The Wolverine (film) X-Men: First Class John Paesano William Stryker

Table 10: Top-5 related concepts from CRX & CCX models for sample target concepts

6. Discussion & Conclusion

In this paper we proposed two models for learning neural embeddings of explicit concepts based on the skip-gram model. Explicit concepts are lexical expressions (single or multi-words) that denote an idea, event, or an object and typically have a set of properties associated with it. In the models presented here, our concept space is the set of all *Wikipedia* article titles. We proposed learning concept representations from concept mentions/references in *Wikipedia* making our models applicable to other open domain and domain specific free-text corpora by firstly wikifying¹³ the text and then learning from concept mentions.

It is clear from the presented results that, the CRX model outperforms the CCX model on tasks that require topical coherence among the concepts vectors (e.g. concept categorization), while the CCX model is advantageous in tasks that require topical relatedness (e.g., measuring entity relatedness). To better show this difference qualitatively, we present a qualitative analysis of both models in Table 10 (target concepts are similar to those reported by Hu et al. (2015)).

As we can notice, the CRX model tends to emphasize concept *topical and categorical similarity*, while the CCX model tends to more emphasize *concept relatedness*. For example, the top-5 concepts closest to "*Harvard University*" using CRX are all universities. While, the CCX model top-5 concepts include, besides educational institutions, location ("*Cambridge, Massachusetts*") and an affiliated group ("*Harvard Society of Fellows*"). The same

13. Wikification is the process of identifying mentions of concepts and entities in a given free-text and linking them to *Wikipedia*

pattern can be noticed with the "X-Men" movie where we get similar genre movies with CRX. While we get related characters such as "William Stryker"¹⁴ with CCX.

Based on these observations, we claim that the CCX model would be beneficial in situations where *diversity* is more desired than *topical coherence*. This claim is also supported by the results we obtained on the concept categorization and dataless densification tasks. On concept categorization, the performance gap between CRX and CCX was large with almost all datasets. On dataless classification, the performance gap was large with documents belonging topics with nuance differences (i.e., *Autos* vs. *Motorcycles*), but with other classes which have clear distinctions, the CCX performance was very competitive to CRX (e.g., *Hockey* vs *Baseball*).

In this paper, we also proposed an *efficient* and *effective* mechanism for BoC densification which outperformed the previously proposed densification schemes on dataless document classification. Unlike these previous densification mechanisms, our method *scales linearly* with the number of the BoC dimensions. In addition, we demonstrated through the results how this efficient mechanism allows generating high quality dense BoC from few concepts alleviating the need of obtaining hundreds of concepts when generating the BoC in the first place.

Our learning method does not require training on a hierarchical concept category graph and is not tightly coupled to linked knowledge base. Rather, we learn concept representations using mentions in free-text corpora with annotated concept mentions which even if not available could be obtained through state-of-the-art entity linking systems.

Finally, the work presented in this paper serves two of our objectives: 1) it demonstrates utilizing textual knowledge bases to learn robust concept embeddings and hence increasing the *effectiveness* of the BoC representation to better capture semantic similarities between textual structures, and 2) it demonstrates utilizing the learned distributed concept vectors to increase the *efficiency* of the semantic representations in terms of space and computational complexities.

References

- Baroni, M., & Lenci, A. (2010). Distributional memory: A general framework for corpus-based semantics. *Computational Linguistics*, 36(4), 673–721.
- Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., & Yakhnenko, O. (2013). Translating embeddings for modeling multi-relational data. In *Advances in neural information processing systems*, pp. 2787–2795.
- Ceccarelli, D., Lucchese, C., Orlando, S., Perego, R., & Trani, S. (2013). Learning relatedness measures for entity linking. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, pp. 139–148. ACM.
- Chang, M.-W., Ratnoff, L.-A., Roth, D., & Srikumar, V. (2008). Importance of semantic representation: Dataless classification.. In *AAAI*, Vol. 2, pp. 830–835.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., & Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug), 2493–2537.

14. https://en.wikipedia.org/wiki/William_Stryker

- Gabrilovich, E., & Markovitch, S. (2007). Computing semantic relatedness using wikipedia-based explicit semantic analysis.. In *IJCAI*, Vol. 7, pp. 1606–1611.
- Harris, Z. S. (1954). Distributional structure.. *Word*.
- Hoffart, J., Seufert, S., Nguyen, D. B., Theobald, M., & Weikum, G. (2012). Kore: keyphrase overlap relatedness for entity disambiguation. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, pp. 545–554. ACM.
- Hu, Z., Huang, P., Deng, Y., Gao, Y., & Xing, E. P. (2015). Entity hierarchy embedding. In *Proceedings of The 53rd Annual Meeting of the Association for Computational Linguistics*.
- Hua, W., Wang, Z., Wang, H., Zheng, K., & Zhou, X. (2015). Short text understanding through lexical-semantic analysis. In *Data Engineering (ICDE), 2015 IEEE 31st International Conference on*, pp. 495–506. IEEE.
- Huang, H., Heck, L., & Ji, H. (2015). Leveraging deep neural networks and knowledge graphs for entity disambiguation. *arXiv preprint arXiv:1504.07678*.
- Hulpus, I., Prangnawarat, N., & Hayes, C. (2015). Path-based semantic relatedness on linked data and its use to word and entity disambiguation. In *International Semantic Web Conference*, pp. 442–457. Springer.
- Kim, D., Wang, H., & Oh, A. H. (2013). Context-dependent conceptualization.. In *IJCAI*, pp. 2330–2336.
- Lang, K. (1995). Newsweeder: Learning to filter netnews. In *Proceedings of the 12th international conference on machine learning*, pp. 331–339.
- Li, P., Wang, H., Zhu, K. Q., Wang, Z., & Wu, X. (2013). Computing term similarity by large probabilistic isa knowledge. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pp. 1401–1410. ACM.
- Li, Y., Zheng, R., Tian, T., Hu, Z., Iyer, R., & Sycara, K. (2016). Joint embedding of hierarchical categories and entities for concept categorization and dataless classification. *arXiv preprint arXiv:1607.07956*.
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013a). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119.
- Milne, D., & Witten, I. H. (2008). Learning to link with wikipedia. In *Proceedings of the 17th ACM conference on Information and knowledge management*, pp. 509–518. ACM.
- Papadimitriou, C. H., & Steiglitz, K. (1982). *Combinatorial optimization: algorithms and complexity*. Courier Corporation.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct), 2825–2830.

- Peng, H., Song, Y., & Roth, D. (2016). Event detection and co-reference with minimal supervision.. EMNLP.
- Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global vectors for word representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12, 1532–1543.
- Rocchio, J. J. (1971). Relevance feedback in information retrieval..
- Schuhmacher, M., & Ponzetto, S. P. (2014). Knowledge-based graph document modeling. In *Proceedings of the 7th ACM international conference on Web search and data mining*, pp. 543–552. ACM.
- Shalaby, W., & Zadrozny, W. (2015). Measuring semantic relatedness using mined semantic analysis. *arXiv preprint arXiv:1512.03465*.
- Song, Y., & Roth, D. (2014). On dataless hierarchical text classification.. In *AAAI*, pp. 1579–1585.
- Song, Y., & Roth, D. (2015). Unsupervised sparse vector densification for short text similarity. In *Proceedings of NAACL*.
- Song, Y., Wang, H., Wang, Z., Li, H., & Chen, W. (2011). Short text conceptualization using a probabilistic knowledgebase. In *Proceedings of the Twenty-Second international joint conference on Artificial Intelligence-Volume Volume Three*, pp. 2330–2336. AAAI Press.
- Song, Y., Wang, S., & Wang, H. (2015). Open domain short text conceptualization: A generative+ descriptive modeling approach.. In *IJCAI*, pp. 3820–3826.
- Wang, Z., & Wang, H. (2016). Understanding short texts. In *the Association for Computational Linguistics (ACL) (Tutorial)*.
- Wang, Z., Wang, H., & Hu, Z. (2014). Head, modifier, and constraint detection in short texts. In *Data Engineering (ICDE), 2014 IEEE 30th International Conference on*, pp. 280–291. IEEE.
- Wang, Z., Zhao, K., Wang, H., Meng, X., & Wen, J.-R. (2015). Query understanding through knowledge-based conceptualization..
- Witten, I., & Milne, D. (2008). An effective, low-cost measure of semantic relatedness obtained from wikipedia links. In *Proceeding of AAAI Workshop on Wikipedia and Artificial Intelligence: an Evolving Synergy*, AAAI Press, Chicago, USA, pp. 25–30.
- Wu, W., Li, H., Wang, H., & Zhu, K. Q. (2012). Probbase: A probabilistic taxonomy for text understanding. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, pp. 481–492. ACM.
- Yamada, I., Shindo, H., Takeda, H., & Takefuji, Y. (2016). Joint learning of the embedding of words and entities for named entity disambiguation. *arXiv preprint arXiv:1601.01343*.
- Zwillinger, D., & Kokoska, S. (1999). *CRC standard probability and statistics tables and formulae*. CRC.