

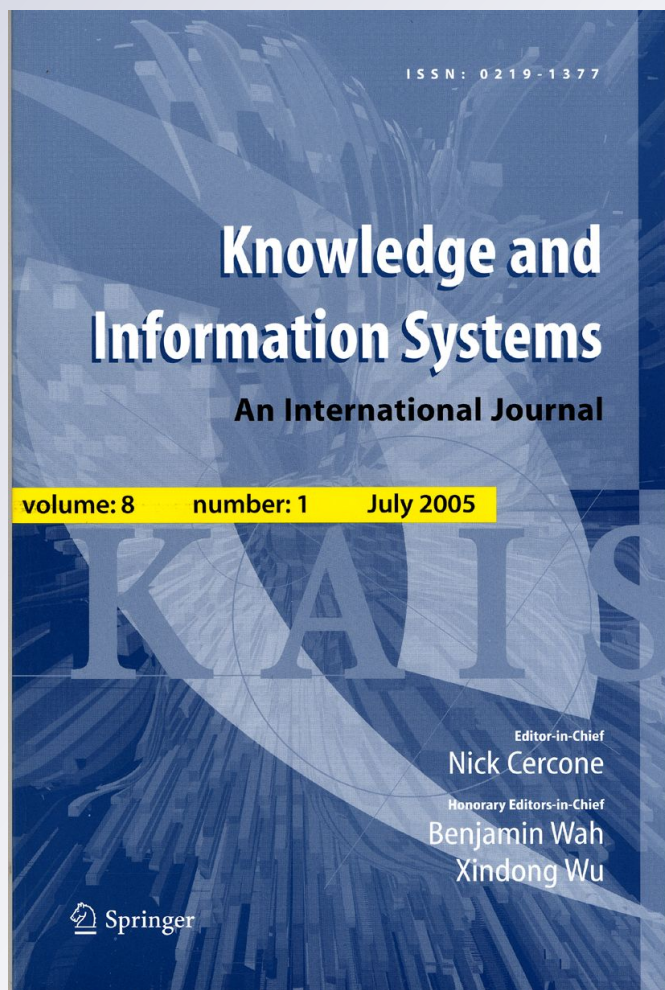
Text document clustering using global term context vectors

Argyris Kalogeratos & Aristidis Likas

Knowledge and Information Systems
An International Journal

ISSN 0219-1377
Volume 31
Number 3

Knowl Inf Syst (2012) 31:455-474
DOI 10.1007/s10115-011-0412-6



Your article is protected by copyright and all rights are held exclusively by Springer-Verlag London Limited. This e-offprint is for personal use only and shall not be self-archived in electronic repositories. If you wish to self-archive your work, please use the accepted author's version for posting to your own website or your institution's repository. You may further deposit the accepted author's version on a funder's repository at a funder's request, provided it is not made publicly available until 12 months after publication.

Text document clustering using global term context vectors

Argyris Kalogeratos · Aristidis Likas

Received: 13 May 2010 / Revised: 20 December 2010 / Accepted: 6 May 2011 /
Published online: 28 May 2011
© Springer-Verlag London Limited 2011

Abstract Despite the advantages of the traditional vector space model (VSM) representation, there are known deficiencies concerning the *term independence* assumption. The high dimensionality and sparsity of the text feature space and phenomena such as polysemy and synonymy can only be handled if a way is provided to measure term similarity. Many approaches have been proposed that map document vectors onto a new feature space where learning algorithms can achieve better solutions. This paper presents the global term context vector-VSM (GTCV-VSM) method for text document representation. It is an extension to VSM that: (i) it captures local contextual information for each term occurrence in the term sequences of documents; (ii) the local contexts for the occurrences of a term are combined to define the global context of that term; (iii) using the global context of all terms a proper semantic matrix is constructed; (iv) this matrix is further used to linearly map traditional VSM (Bag of Words—BOW) document vectors onto a ‘*semantically smoothed*’ feature space where problems such as text document clustering can be solved more efficiently. We present an experimental study demonstrating the improvement of clustering results when the proposed GTCV-VSM representation is used compared with traditional VSM-based approaches.

Keywords Text mining · Document clustering · Semantic matrix · Data projection

1 Introduction

The text document clustering procedure aims toward automatically partitioning a given collection of unlabeled text documents into a (usually predefined) number of groups, called

A. Kalogeratos · A. Likas (✉)
Department of Computer Science, University of Ioannina,
45110, Ioannina, Greece
e-mail: arly@cs.uoi.gr

A. Kalogeratos
e-mail: akaloger@cs.uoi.gr

clusters, such that similar documents are assigned to the same cluster while dissimilar documents are assigned to different clusters. This is a task that discovers the underlying structure in a set of data objects and enables the efficient organization and navigation in large text collections.

The challenging characteristics of the text document clustering problem are related to the complexity of the natural language. Text documents are represented in high dimensional and sparse (*HDS*) feature spaces, due to their large term vocabularies (the number of different terms of a document collection or text features in general). In an HDS feature space, the difference between the distance of two similar objects and the distance of two dissimilar objects is relatively small [4]. This phenomenon prevents clustering methods from achieving good data partitions. Moreover, the text semantics, e.g., term correlations, are mostly implicit and non-trivial, hence difficult to extract without prior knowledge for a specific problem.

The traditional document representation is the *vector space model (VSM)* [28] where each document is represented by a vector of weights corresponding to text features. Many variations of VSM have been proposed [17] that differ in what they consider as a feature or '*term*'. The most common approach is to consider different words as distinct terms, which is the widely known *Bag Of Words (BOW)* model. An extension is the *Bag Of Phrases (BOP)* model [23] that extracts a set of informative phrases or *word n-grams* (n consecutive words). Especially for noisy document collections, e.g., containing many spelling errors, or collections whose language is not known in advance, it is often better to use VSM to model the distribution of *character n-grams* in documents. Herein, we consider word features and we refer to them as *terms*; however, the procedures we describe can be directly extended to more complex features.

Despite the simplicity of the popular word-based VSM version, there are common language phenomena that it cannot handle. More specifically, it cannot distinguish the different senses of a polysemous word in different contexts or realize the common sense between synonyms. It also fails to recognize multi-word expressions (e.g., '*Olympic Games*'). These deficiencies are in part due to the over-simplistic assumption of *term independence*, where each dimension of the HDS feature space is considered to be vertical to the others and makes the classic VSM model incapable of capturing the complex language semantics. The VSM representations of documents can be improved by examining the relations between terms either at a low level, such as *terms co-occurrence frequency*, or at a higher *semantic similarity* level.

Among the popular approaches is the *Latent Semantic Indexing (LSI)* [7] that solves an eigenproblem using *Singular Value Decomposition (SVD)* to determine a proper feature space to project data. *Concept Indexing* [16] computes a k -partition by clustering the documents and then uses the centroid vectors of the clusters as the axes of the reduced space. Similarly, *Concept Decomposition* [8] approximates in a least-squares fashion the term-by-document data matrix using centroid vectors. A more simple but quite efficient method is the *Generalized vector space model (GVSM)* [31]. GSVM represents documents in the document similarity space, i.e., each document is represented as a vector containing its similarities to the rest of the documents in the collection. The *Context Vector Model (CVM-VSM)* [5] is a VSM-extension that describes the semantics of each term by introducing a *term context vector* that stores its similarities to the other terms. The similarity between terms is based on a document-wise term co-occurrence frequency. The term context vectors are then used to map document vectors into a feature space of equal size to the original, but less sparse. The *ontology-based VSM* approaches [13, 15] map the terms of the original space onto a feature space defined by a hierarchically structured thesaurus, called *ontology*. Ontologies provide information about the words of a language and their possible semantic relations; thus, an efficient mapping can

disambiguate the word senses in the context of each document. The main disadvantage is that, in most cases, the ontologies are static and rather generic knowledge bases, which may cause heavy semantic smoothing of the data. A special text representation problem is related to very short texts [14, 25].

In this work, we present the *Global Term Context Vector-VSM (GTCV-VSM)* representation that is an entirely corpus-based extension to the traditional VSM that incorporates *contextual information* for each vocabulary term. First, the *local context* for each term occurrence in the term sequences of documents is captured and represented in vector space by exploiting the idea of the *Locally Weighted Bag of Words* [18]. Then, all the local contexts of a term are combined to form its *global context vector*. Global context vectors constitute a *semantic matrix* that efficiently maps the traditional VSM document vectors onto a semantically richer feature space of same dimensionality to the original. As indicated by our experimental study, in the new space, superior clustering solutions are achieved using well-known clustering algorithms such as the spherical k -means [8] or spectral clustering [24].

The rest of this paper is organized as follows. Section 2 provides some background on document representation using the vector space model. In Sect. 3, we describe recent approaches for representing a text document using histograms that describe the local context at each location of the document-term sequence. In Sect. 4, we present our proposed approach for document representation. The experimental results are presented in Sect. 5, and finally, in Sect. 6, we provide conclusions and directions for future work.

2 Document representation in vector space

In order to apply any clustering algorithm, the raw collection of N text documents must be first preprocessed and represented in a suitable feature space. A standard approach is to eliminate trivial words (e.g., stopwords) and words that appear in a small number of documents. Then, *stemming* [26] is applied, which aims to replace each word by its corresponding word *stem*. The V derived word stems constitute the collection's *term vocabulary*, denoted as $\mathcal{V} = \{v_1, \dots, v_V\}$. Thus, a text document, which is a finite term sequence of T vocabulary terms, is denoted as $d^{seq} = \langle d^{seq}(1), \dots, d^{seq}(T) \rangle$, with $d^{seq}(i) \in \mathcal{V}$. For example, the phrase 'The dog ate a cat and a mouse!' is a sequence $d^{seq} = \langle \text{dog}, \text{ate}, \text{cat}, \text{mouse} \rangle$.

2.1 The bag of words model

According to the typical VSM approach, the *Bag of Words (BOW)* model, a document is represented by a vector $d \in \mathbb{R}^V$, where each word term v_i of the vocabulary is associated with a single vector dimension. The most popular weighting scheme is the *normalized $tf \times idf$* that introduces the *inverse document frequency* as an external weight to enforce the terms that have *discrimination power* and appear in a small number of documents. For the v_i vocabulary term, it is computed as $idf_i = \log(N/df_i)$, where N denotes the total number of documents and df_i denotes the *document frequency*, i.e., the number of documents that contain term v_i . Thus, the normalized $tf \times idf$ BOW vector is a mapping of the term sequence d^{seq} defined as follows

$$\Phi_{bow} : d^{seq} \rightarrow d = h \cdot (tf_1 idf_1, \dots, tf_V idf_V)^\top \in \mathbb{R}^V, \quad (1)$$

where normalization is performed with respect to the Euclidean norm using the coefficient h . The document collection can then be represented using the N document vectors as rows in

the *Document-Term matrix* D , which is a $N \times V$ matrix whose rows and columns are indexed by the documents and the vocabulary terms, respectively.

In the VSM, there are several alternatives to quantify the semantic similarity between document pairs. Among them, *Cosine similarity* has shown to be an effective measure [11], and for a pair of document vectors, d_i and d_j is given by

$$sim_{cos}(d_i, d_j) = \frac{d_i^\top d_j}{\|d_i\|_2 \|d_j\|_2} \in [0, 1]. \tag{2}$$

Unit similarity value implies that the two documents are described by identical distributions of term frequencies. Note that this is equal to the dot product $d_i^\top d_j$ if document vectors are normalized in the unit positive V -dimensional hypersphere.

2.2 Extensions to VSM

The BOW model, despite having a series of advantages, such as generality and simplicity, it cannot model efficiently the rich semantic content of text. The Bag Of Phrases model uses phrases of two or three consecutive words as features. Its disadvantage is the fact that it has been observed that as phrases become longer, they obtain superior semantic value, but at the same time, they become statistically inferior with respect to single-word representations [19]. A category of methods developed aiming on tackling this difficulty recognize the frequent *wordsets* (unordered itemsets) in a document collection [3, 10], while the method proposed in the study by [20] exploits the frequent word subsequences (ordered) that are stored in a *Generalized Suffix Tree (GST)* for each document.

Modern variations of VSM are used to tackle the difficulties occurring due to HDS spaces, by projecting the document vectors onto a new feature space called *concept space*. Each *concept* is represented as a *concept vector* of relations between the concept and the vocabulary terms. Generally, this approach of document mapping can be expressed as

$$\Phi_{VSM} : d \rightarrow d' = Sd \in \mathbb{R}^{V'}, \quad V' \leq V, \tag{3}$$

where the $V' \times V$ matrix S stores the concept vectors as rows. This projection matrix is also known as *semantic matrix*. The Cosine similarity between two normalized document images in the concept space can be computed as a dot product

$$sim_{sem}^{(cos)}(d'_i, d'_j) = (\overline{Sd_i})^\top (\overline{Sd_j}) = \left(h_i^S Sd_i\right)^\top \left(h_j^S Sd_j\right) = h_i^S h_j^S \left(d_i^\top S^\top Sd_j\right), \tag{4}$$

where the scalar normalization coefficient for each document is $h_i^S = 1/\|Sd_i\|_2$. The similarity defined in Eq. 4 can be interpreted in two ways: (i) as a dot product of the document images $(\overline{Sd_i})^\top (\overline{Sd_j})$ that both belong to the new space $\mathbb{R}^{V'}$ and (ii) as a composite measure that takes into account the pairwise correlations between the original features expressed by the matrix $S^\top S$.

There is a variety of methods proposing alternative ways to define the semantic matrix though many of them are based on the above linear mapping. The widely used *Latent Semantic Indexing (LSI)* [7] projects the document vectors onto a space spanned by the eigenvectors corresponding to the V' largest eigenvalues of the matrix $D^\top D$. The eigenvectors are extracted by the means of *Singular Value Decomposition (SVD)* on matrix D^\top , and they capture the latent semantic information of the feature space. In this case, each eigenvector is a different concept vector and V' is a user parameter much smaller than V , while there is also a considerable computational cost to perform the SVD. In *Concept Indexing* [16], the concept

vectors are the centroids of a V' -partition obtained by applying document clustering. In [9], statistical information such as the covariance matrix is combined with traditional mapping approaches into latent space (LSI, PCA) to compose a hybrid vector mapping.

A computationally simpler alternative that utilizes the Document-Term Matrix D as a semantic matrix is the *Generalized vector space model (GVSM)* [31], i.e., $S_{gvsm} = D$ and the image of a document is given by $d' = Dd$. By examining the product $Dd \in \mathbb{R}^{N \times 1}$, we can conclude that a GVSM projected document vector d' has lower dimensionality if $N \leq V$. Moreover, if both d and D are properly normalized, then image vector d' consists of the N Cosine similarities between the document vector d and the rest of the $N - 1$ documents in the collection. This observation implies that the GVSM works in the *document similarity space* by considering each document as a different concept. On the other hand, the respective product $S_{gvsm}^T S_{gvsm} = D^T D$ (used in Eq. 4) is a $V \times V$ *Term Similarity Matrix* whose r -th row has the dot-product similarities between term v_r and the rest of the $V - 1$ of vocabulary terms. Note that terms become more similar as their corresponding normalized frequency distributions into the N documents are more alike. Based on the GVSM model, it is proposed in [1] to build local semantic matrices for each cluster during document clustering.

A rather different approach proposed in [5] for information retrieval is the *Context Vector Model (CVM-VSM)* where, instead of a few concise concept vectors, it computes the context in which each of the V vocabulary terms appears in the data set, called *term context vector (tcv)*. This model computes a $V \times V$ matrix S_{cvm} containing the term context vectors as rows. Each tcv_i vector aims to capture the V pairwise similarities of term v_i to the rest of the vocabulary terms. Such similarity is computed using a *co-occurrence frequency measure*. Each matrix element $[S_{cvm}]_{ij}$ stores the similarity between terms v_i and v_j computed as

$$[S_{cvm}]_{ij} = \begin{cases} 1 & , \quad i = j \\ \frac{\sum_{r=1}^N tf_{ri} tf_{rj}}{\sum_{r=1}^N (tf_{ri} \cdot \sum_{q=1, q \neq i}^V tf_{rq})} & , \quad i \neq j. \end{cases} \quad (5)$$

Note that this measure is not symmetric, generally $[S_{cvm}]_{ij} \neq [S_{cvm}]_{ji}$, due to the denominator that normalizes the pairwise similarity to $[0, 1]$ with respect to the 'total amount' of similarity between term v_i and the other vocabulary terms. The rows of matrix S_{cvm} can be normalized with respect to the Euclidean norm, and each document image is then computed as the centroid of the normalized context vectors of all terms appearing in that document

$$\Phi_{cvm} : d \rightarrow d' = \sum_{i=1}^V tf_i \cdot tcv_i, \quad (6)$$

where tf_i is the frequency of term v_i . The motivation for using term context vectors is to capture the semantic content of a document based on the co-occurrence frequency of terms in the same document, averaged over the whole corpus. The CVM-VSM representation is less sparse than BOW. Moreover, weights such as *idf* can be incorporated to the transformed document vectors computed using Eq. 6. In [5], several more complicated weighting alternatives have been tested in the context of information retrieval that in our text document, clustering experiments did not perform better than the standard *idf* weights.

In a higher semantic level than term co-occurrences, additional information for vocabulary terms provided by ontologies has also been exploited to compute the term similarities and to construct a proper semantic matrix. *WordNet* [22] and *Wikipedia* [30] have been used for this purpose in [6, 15], and [29], respectively.

2.3 Discussion

Summarizing the properties of the above-mentioned vector-based document representations, in the traditional BOW approach, the dimensions of the term feature space are considered to be independent to each other. Such an assumption is very simplistic, since there exist semantic relations among terms that are ignored. The VSM-extensions aim to achieve *semantic smoothing*, a process that redistributes the term weights of a vector model, or map data in a new feature space, by taking into account the correlations between terms. For instance, if the term 'child' appears in a document, then it could be assumed that the term 'kid' is also related to the specific document or even terms like 'boy', 'girl', 'toy'. The resulting representation model is also a VSM, but the document vectors become less sparse and the independence of features is mitigated in an indirect way. The smoothing is usually achieved by a linear mapping of data vectors to a new feature space using a semantic matrix S . It is convenient to think that the new document vector $d' = Sd$ contains the dot product similarities between the original BOW vector d and the rows of the semantic matrix S .

A basic difference between the various semantic smoothing methods is related to the dimension of the new feature space, which is determined by the number V' of row vectors of matrix S . In case their number is less than the size V of the vocabulary, such vectors are called as *concept vectors* and are usually produced using the LSI method. Each concept vector has a distribution of weights associated with the V original terms that define their contribution of to the corresponding *concept*. Of course, the resulting representation of the smoothed vector d' is less interpretable than the original, and there is always a problem of determining the proper number of concept vectors.

An alternative approach for semantic smoothing assumes that each row vector of matrix S is associated with one vocabulary term. Unlike a concept vector that describes abstract semantics of higher level, here, the elements of each vector describe the relation of this term to the other terms. Those relations constitute the so-called *term context*, thus the respective vector is called *term context vector*. Each element of the mapped vector d' will contain the dot product similarity between document d and the corresponding term context vector, i.e., for each term v_i , the element d'_i provides the degree to which the original document d contains the term v_i and its context, instead of just computing its frequency as happens in the BOW representation. Note also that in BOW representation, a dot product would give zero similarity for two documents that do not have common terms. On the contrary, the dot product between a document vector and a term context vector of a term v_i that does not appear in that document may give a non-zero similarity. This happens if the document contains at least one term v_j with non-zero weight in the context of term v_i . For this reason, the smoothed representation d' is usually less sparse than d and retains their interpretability of dimensions. Moreover, concept-based methods may be applied on the new representations.

The motivation of our work is to establish the importance of *term context vectors* and to define an efficient way to compute them. The CVM-SVM method considers that the term context is computed based on term co-occurrence frequency at the document level. It does not take into account the sequential nature of text and thus ignores the local distance of terms when computing term context. On the other hand, the GTCV-VSM proposed in this work extends the previous approach by considering term context at three levels: (i) It uses the notion of *local term context vector (ltcv)* to model the context around the location in the text sequence where a term appears. These vectors are computed using a local smoothing kernel as suggested in the *LoWBOW* approach [18] which is described in the next section. The kernel *takes into account the distance* in which other terms appear around the sequence location under consideration. (ii) It computes

the *document term context vector* (d_{tcv}) for each term that summarizes the term context at the document level, and (iii) it computes the final *global term context vector* (g_{tcv}) for each term representing the overall term context at corpus level. The g_{tcv} vectors constitute the rows of the semantic matrix S . Thus, the intuition behind GTCV-VSM approach is to capture the local term context from term sequences and then to construct a representation for global term context by averaging l_{tcvs} at the document and corpus level.

3 Utilizing local contextual information

A text document can be considered as a finite term sequence of its T consecutive terms denoted as $d^{seq} = \langle d^{seq}(1), \dots, d^{seq}(T) \rangle$ but, except for Bag of Phrases, so far in this paper, the previously mentioned VSM-extensions ignore this property. A category of methods has been proposed aiming to capture local information directly from the term sequence of a document. The representation proposed that in the study by [27], first considers a segmentation of the sequence that is done by dragging a window of n terms along the sequence and computing the local BOW vectors for each of the overlapping segments. All these local BOW vectors constitute the document representation called *Local Word Bag* (LWB). To compute the similarity between a pair of documents, the authors introduce a variant of the *VG-Pyramid Matching Kernel* [12] that maps the two sets of local BOW vectors to a multi-resolution histogram and computes a weighted histogram intersection.

Another approach for text representation presented in [18] is the *Locally Weighted Bag of Words* ($LoWBOW$) that preserves local contextual information of text documents by the effective modeling of the text sequential structure. At first, a number of L equally distant locations are defined in the term sequence. Each sequence location $\ell_i, i = 1, \dots, L$ is then associated with a local histogram which is a point in the multinomial simplex \mathbb{P}_{V-1} , where V is the number of vocabulary terms. More specifically, for $(V - 1) \geq 0$, the \mathbb{P}_{V-1} space is the $(V - 1)$ -dimensional subset of \mathbb{R}^V that contains all probability vectors (histograms) over V objects (for a discussion on the multinomial simplex see the Appendix of [18])

$$\mathbb{P}_{V-1} = \left\{ H \in \mathbb{R}^V : H_i \geq 0, \forall i = 1, \dots, V \text{ and } \sum_{i=1}^V H_i = 1 \right\}. \tag{7}$$

Contrary to LWB , in $LoWBOW$ the local histogram is computed using a *smoothing kernel* to weight the contribution of terms appearing around the referenced location in the term sequence and to assign more importance to closely neighboring terms. Denoting as $H_{\delta(d^{seq}(t))}$ the *trivial term histogram* of V terms whose probability mass is concentrated only at the term that occurs at the location t in d^{seq}

$$[H_{\delta(d^{seq}(t))}]_i = \begin{cases} 1, & v_i = d^{seq}(t) \\ 0, & v_i \neq d^{seq}(t) \end{cases}, \quad i = 1, \dots, V, \tag{8}$$

then the locally smoothed histogram at a location ℓ in the d^{seq} term sequence is computed as in [18]

$$lowbow(d^{seq}, \ell) = \sum_{t=1}^T H_{\delta(d^{seq}(t))} K_{\ell, \sigma}(t), \tag{9}$$

where T is the length of d^{seq} . $K_{\ell,\sigma}(t)$ denotes the weight for location t in sequence given by a discrete Gaussian weighting kernel function of mean value ℓ and standard deviation σ . Specifically, the weighting function is a Gaussian probability density function restricted in $[1, T]$ and renormalized so that $\sum_{t=1}^T K_{\ell,\sigma}(t) = 1$. It is easy to verify that the result of the histogram smoothing of Eq. 9 is also a histogram.

It must be noted that for $\sigma = 0$, the *lowbow* histogram (Eq. 9) coincides with the trivial histogram $H_{\delta(d^{seq}(\ell))}$, where all the probability mass is concentrated at the term at location ℓ . As σ grows, part of the probability mass is transferred to the terms occurring near location ℓ . In this way, the *lowbow* histogram at location ℓ is enriched with information about the terms occurring in the neighborhood of ℓ . The smoothing parameter σ adjusts the ‘locality’ of term semantics that is taken into account by the model. Thus, instead of mining unordered local vectors as in [27], the LoWBOW approach embeds the term sequence of a document in the \mathbb{P}_{V-1} simplex. The sequence of the L locally smoothed histograms (denoted as *lowbow histograms*) form a curve in the $(V - 1)$ -dimensional simplex (denoted as *LoWBOW curve*). Figure 1 illustrates the LoWBOW curves generated for a toy example and describes the role of parameter σ . In this figure, we aim to illustrate (i) the LoWBOW curve representation, i.e., the curve that corresponds to a sequence of histograms (local context vectors), where each local context vector is computed at a specific location of the sequence and corresponds to a point in the $(V - 1)$ -dimensional simplex; (ii) the impact of the smoothing coefficient σ on the computed local context vectors. It is illustrated that the increase in smoothing makes the *lowbow* histograms (points of the curve) more similar. This can also be verified by observing that as smoothing increases, the curve becomes more concentrated around a central location of the simplex. For $\sigma = \infty$, all histograms become similar to the BOW representation and the curve reduces to a single point. On the contrary, for $\sigma = 0$, the histograms correspond to simplex corners.

A similarity measure between LoWBOW curves has been proposed in [18] that assumes a sequential correspondence between two documents and computes the sum of the similarities between the L pairs of LoWBOW histograms. Obviously, it is expected for this similarity measure to underestimate the thematic similarity between documents that follow different order in the presentation of similar semantic content.

4 A semantic matrix based on global term context vectors

In this section, we present the global term context vector-VSM (GTCV-VSM) approach for capturing the semantics of the original term feature space of a document collection. The method computes the contextual information of each vocabulary term, which is subsequently utilized in order to create a semantic matrix. In analogy with CVM-VSM, our approach reduces data sparsity but not dimensionality. The interpretability of the derived vector dimensions remains as strong as in the BOW model as the value of each dimension of the mapped vector corresponds to one vocabulary term. Methods that reduce data dimensionality could also be applied on the new representations at a subsequent phase. Compared with CVM-VSM, GTCV-VSM generalizes the way the term context is computed by taking into account the distance between terms in the term sequence of each document. This is achieved by exploiting the idea of LoWBOW to describe the local contextual information at a certain location in a term sequence. It must be noted that our method borrows from the LoWBOW approach only the way the local histogram is computed at each location of the term sequence and does not make use of the LoWBOW curve representation.

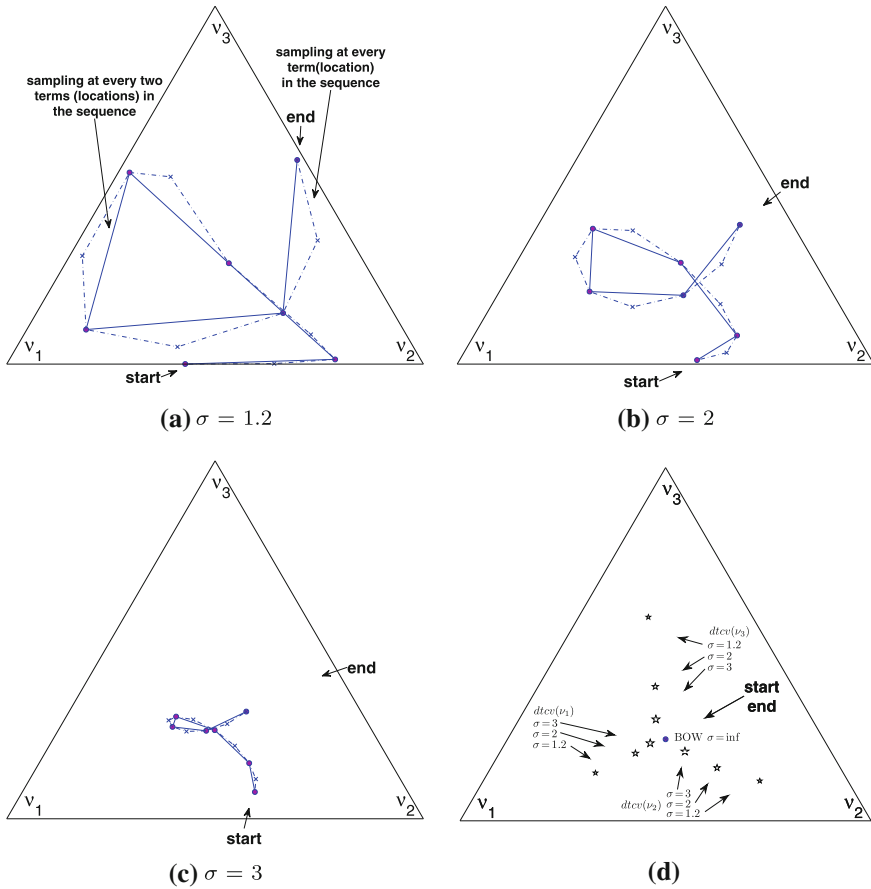


Fig. 1 A toy example where the sequence $\langle v_1, v_2, v_2, v_2, v_1, v_3, v_3, v_1, v_1, v_1, v_2, v_2, v_3 \rangle$ is considered that uses three different terms v_1, v_2, v_3 (vocabulary size: $V = 3$). The subfigures present *LoWBOW* curves in the $(V - 1)$ -dimensional simplex for increasing values of the parameter σ that induces more smoothing to the curve. Each point of the curve corresponds to a local histogram computed at a sequence location. The more a term affects the local context at a location in the sequence, the more the curve point (the *lowbow* histogram related to that location) moves toward the respective corner of the simplex. For $\sigma = 0$, local histograms correspond to simplex corners; thus, the curve moves from corner to corner of the simplex. Two different sampling rates for *LoWBOW* representation are illustrated: sampling at every term location in the sequence (*dashed line*), which is the our strategy to collect contextual information for each term, and sampling every two terms (*solid line*). **d** For $\sigma = \infty$, the *LoWBOW* curve reduces to a single point that coincides with the *BOW* histogram of the sequence. In **d**, we present as ‘stars’ the average *ltcv* histograms for each term (*dtcv* histograms) for the three different values of σ and $\alpha = 0.6$ for all terms. As the value of σ increases, the *dtcv* histograms of all terms become more similar tending to coincide with the *BOW* representation

More specifically, we define the *local term context vector* (*ltcv*) as a histogram associated with the exact occurrence of term $d^{seq}(\ell)$ at location ℓ in a sequence d^{seq} . Hence, one *ltcv* vector is computed at every location in the term sequence, i.e., $\ell = 1, \dots, T$. Note that *GTCV-VSM* does not preserve any curve representation. This means that we are not interested in the temporal order of the local term context vectors. The $ltcv(d^{seq}, \ell)$ is a modified *lowbow*(d^{seq}, ℓ) probability vector that represents contextual information around location ℓ , while adjusting explicitly the *self-weight* $\alpha_{d^{seq}(\ell)}$ of the reference term appearing

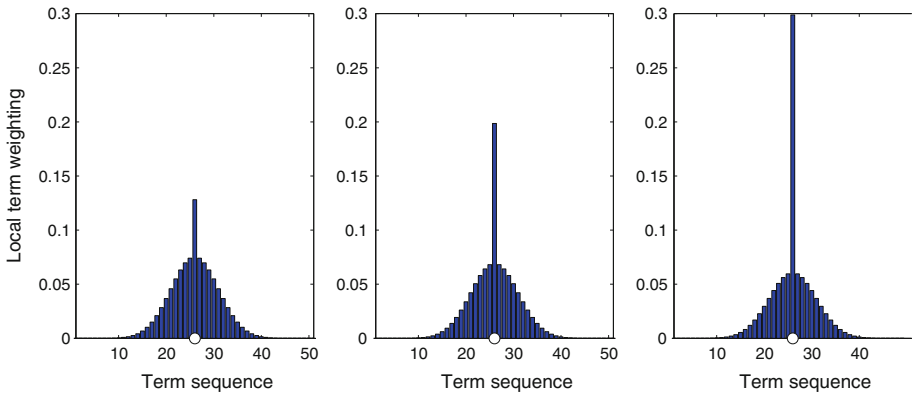


Fig. 2 Various weight distributions for the neighboring terms around a reference term occurring in the middle of a term sequence of length 50. The distributions are obtained by varying the value of parameter α in Eq. 10. This distribution defines the contribution of each term to the context of the specific reference term. The scale value of the local kernel is set to $\sigma = 5$, while self-weight is set to α to 0.05 (left), 0.10 (middle), 0.2 (right)

at location ℓ

$$[l_{tcv}(d^{seq}, \ell)]_i = \begin{cases} \alpha_{d^{seq}(\ell)} & , \quad v_i = d^{seq}(\ell), \\ (1 - \alpha_{d^{seq}(\ell)}) \cdot \frac{idf_i \cdot [lowbow(d^{seq}, \ell)]_i}{\sum_{j=1, j \neq i}^V idf_j \cdot [lowbow(d^{seq}, \ell)]_j} & , \quad v_i \neq d^{seq}(\ell). \end{cases} \quad (10)$$

The self-weight ($0 \leq \alpha_{d^{seq}(\ell)} \leq 1$) adjusts the relative importance between contextual information (computed using the *lowbow* histogram) and the self-representation of each term. Figure 2 illustrates an example of how the value of parameter α affects the local term weighting around a reference term in a sequence. When the parameter σ of the Gaussian smoothing kernel is set to zero, or $\alpha = 1$, the $l_{tcv}(d^{seq}, \ell)$ reduces to a trivial histogram $H_{\delta(d^{seq}(\ell))}$ (see Eq. 8). The other extreme is the infinite σ value, where for small α values, all the l_{tcv} computed in a document d become similar to the *tf* histogram for that document.

The latter observation is the reason for considering an explicit self-weight in Eq. 10, because a flat smoothing kernel obtained for large σ value can make a *lowbow* vector to have improperly low self-weight for the reference term. For example, if a term appears once in a document, then the *lowbow* vector with $\sigma = \infty$ at that location would contain very low weight for that term. Generally, the value of α_v determines how much the context vector of term v should be dominated by the self-weight of term v . In our method, we set this parameter independently for each individual term as a function of its *idf_v* component

$$\alpha_v = \lambda + (1 - \lambda) \cdot \left(1 - \frac{idf_v}{\log N}\right), \quad \lambda \in [0, 1], \quad (11)$$

where λ is a lower bound for all α_v , $v = 1, \dots, V$ (in our experiments we used $\lambda = 0.2$). The rationale for the above equation is that for terms with high document frequency (i.e., low *idf_v*), we assign high α_v values that suppress the local context in the respective context vectors. In other words, the context is considered more important for terms that occur in fewer documents. In Fig. 3a, we present an example illustrating the l_{tcv} vectors of two term sequences presented in Fig. 3c.

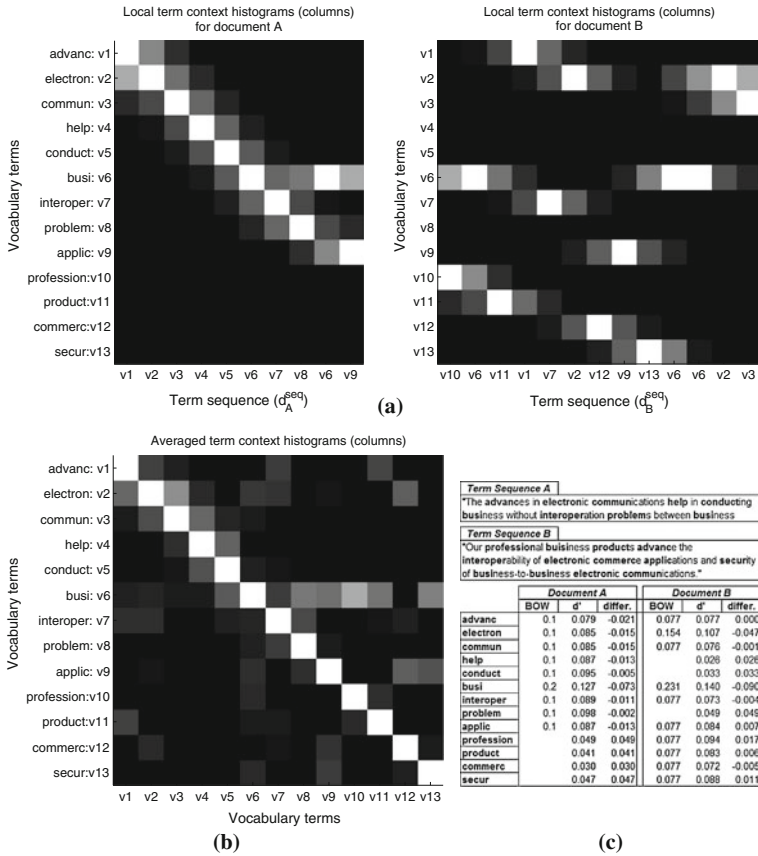


Fig. 3 An example of how *l_tc_v* histograms are used to summarize the overall context in which a term appears in the two term sequences of (c) using Eq. 14. **a** The term sequences (*x*-axis) of documents A, B are presented and the corresponding local term context vectors are illustrated as gray-scaled *columns*. Those vectors are computed at every location in the sequence using a Gaussian smoothing kernel with $\sigma = 1$ and $\alpha = 0.6$ for all terms. Brighter intensity at cell *i, j* indicates higher contribution of the term v_i to the local context of the term appearing at location *j* in the sequence. **b** The resulting transposed semantic matrix (S^T), where the *gray-scaled* columns illustrate the global contextual information for each vocabulary term computed by averaging the respective local context histograms (Eq. 13). **c** The two initial term sequences (the stem of each non-trivial term is emphasized). Assuming the same idf weight for each vocabulary term, the table presents the BOW vector, the transformed vector d' using Eq. 14 as well as the effect of semantic smoothing ($differ = BOW - d'$) on document vectors. The redistribution of term weights that results by the proposed mapping reveals is done in such a way that low-frequency terms are gaining weight against the more frequent ones. Note also that the similarity between the two documents is 0.756 for the BOW model and 0.896 for the GTCV-VSM

We further define the *document term context vector* (*d_tc_v*) as a probability vector that summarizes the context of a specific term at the document level by averaging the *l_tc_v* histograms corresponding to the occurrences of this term in the document. More specifically, suppose that a term *v* appears $no_{i,v} > 0$ times in the term sequence d_i^{seq} (i.e., in the *i*-th document) which is of length T_i . Then, the *d_tc_v* of this term *v* for document *i* is computed as:

$$d_{tcv}(d_i^{seq}, v) = \frac{1}{no_{v,i}} \sum_{j=1}^{no_{i,v}} l_{tcv}(d_i^{seq}, \ell_{i,v}(j)), \quad (12)$$

where $\ell_{i,v}(j)$ is an integer value in $[1, \dots, T_i]$ denoting the location of the j -th occurrence of v in d_i^{seq} .

Next, the *global term context vector* ($gtcv$) is defined for a vocabulary term v so as to represent the overall contextual information for all appearances of v in the corpus of all N term sequences (documents).

$$gtcv(v) = h_{gtcv(v)} \left(\sum_{i=1}^N tf_{i,v} dtcv(d_i^{seq}, v) \right). \tag{13}$$

The coefficient $h_{gtcv(v)}$ normalizes the vector $gtcv(v)$ with respect to the Euclidean norm, and $tf_{i,v}$ is the frequency of the term v in the i -th document. Thus, the $gtcv(v)$ of term v is computed using a weighted average of the document context vectors $dtcv(d_i^{seq}, v)$ obtained for each document i in which term v appears. Thus, in contrast to LoWBOW curve approach which focuses on the sequence of local histograms that describe the writing structure of a document, our method focuses on the extraction of the global semantic context of a term by averaging the local contextual information at all the corpus locations where this term appears.

Finally, the extracted global contextual information is used to construct the $V \times V$ semantic matrix S_{gtcv} where each row v is the $gtcv(v)$ vector of the corresponding vocabulary term v . Figure 1d provides an example of illustrating the $dtcv(d_i^{seq}, v)$ vectors for each document (the points denoted as 'stars'). Figure 3b illustrates the final $gtcv$ vectors obtained by averaging the document level contexts for each vocabulary term.

To map a document using the proposed global term context vector-VSM approach, we compute the vector d' where each element v is Cosine similarity between the BOW representation d of the document and the global term context vector $gtcv(v)$:

$$\Phi_{gtcv} : d \rightarrow d' = S_{gtcv} d, d' \in \mathbb{R}^V. \tag{14}$$

Note that the transformed document vector d' is V -dimensional that retains the interpretability, since each dimension still corresponds to a unique vocabulary term. Moreover, if $\sigma = 0$ and $\alpha > 0$, then $S_{gtcv} d = d$. Looking at Eq. 4, the product $S_{gtcv}^T S_{gtcv}$ essentially computes a Term Similarity Matrix where the similarity between two terms is based on the distribution of term weights in their respective global term context vectors, i.e., on the similarity of their global context histograms. The table of Fig. 3c illustrates the effect of redistribution (compared with BOW) of the term weights (semantic smoothing) in the transformed document vectors achieved by the proposed mapping.

The procedure of representing the input documents using GTCV-VSM takes place in the preprocessing phase. Let T_i the length of the i -th document and V_i its vocabulary. Let also V the size of the whole corpus vocabulary. Then, the cost to compute one l_{tcv} vector at a location of the term sequence using Eq. 10, and to add its V_i non-zero dimensions to the respective d_{tcv} , is $O(T_i + V_i)$. This is done T_i times and the final d_{tcv} of each different term of the document is added to the respective the $gtcv$ rows. Thus, using proper notation for the average length \bar{T}_i and vocabulary size \bar{V}_i of the documents in a corpus, the cost of constructing the semantic matrix can be expressed as $O(N \cdot \bar{T}_i \cdot (\bar{T}_i + 2 \cdot \bar{V}_i))$. However, since $\bar{V}_i \leq \bar{T}_i \ll V$, the overall computational cost of the GTCV-VSM is determined by the $O(N \cdot V^2)$ cost of the matrix multiplication of the mapping of Eq. 14.

Table 1 Characteristics of text document collections

Name	Topics	Classes	N	Balance	V	\overline{V}_i	\overline{T}_i
D_1	20-NGs: graphics, windows.x, motor, baseball, space, mideast	6	2,000	200/400	4,343	48.8	110
D_2	20-NGs: atheism, autos, baseball, electronics, med, mac, motor, politics.misc	7	3,500	500/500	6,442	52.6	108
D_3	20-NGs: atheism, christian, guns, mideast	4	1,600	400/400	4,080	62	131
D_4	20-NGs: forsale, autos, baseball, motor, hockey	5	1,250	250/250	4,762	44.1	104
D_5	Reuters-21,578: acq, corn, crude, earn, grain, interest, money-fx, ship, trade, wheat	10	9979	237/3,964	5,613	39.1	76

N denotes the number of documents, V is the size of the global vocabulary and \overline{V}_i the average document vocabulary, *Balance* is the ratio of the smallest to the largest class and \overline{T}_i is the average length of the term sequences of documents

5 Clustering experiments

Our experimental setup was based on five different data sets: D_1 – D_4 are subsets of the 20-Newsgrroups,¹ while D_5 is the Mod Apte split [2] version of the Reuters-21578² benchmark document collection where the 10 classes with larger number of training examples are kept. The characteristics of these data sets are presented in Table 1. The preprocessing of data sets included the removal of all tags, headers, and metadata from the documents, while applied word stemming and discarded terms appearing in less than five documents. It is worth mentioning how we preprocessed the term sequences of documents. We considered a *dummy term* that replaced in the sequences all the low-frequency terms that were discarded so as to maintain the relative distance between the terms that remained in each sequence. For similar reasons, two dummy terms were considered at the end of every sentence denoted by characters as (e.g., ‘.’, ‘?’ , ‘!’). The dummy term is ignored when constructing the final data vectors.

For each data set, we have considered several data mappings, Φ and after each mapping, the *spherical k-means (spk-means)* [8] and spectral clustering (*spectral-c*) [24] algorithms were applied to cluster the mapped documents vectors into the k predefined number of clusters corresponding to the different topics (classes) in a collection. In contrast to *k-means* that is based on the Euclidean distance [21], spk-means uses the Cosine similarity and maximizes the *Cohesion* of the clusters $C = \{c_1, \dots, c_k\}$

$$Cohesion(C) = \sum_{j=1}^k \sum_{d_i \in c_j} u_j^\top d_i, \tag{15}$$

where u_j is the normalized centroid of cluster c_j with respect to the Euclidean norm.

Spectral clustering projects the document vectors in a subspace that is spanned by the k largest eigenvectors of the Laplacian matrix L computed from the similarity matrix $A^{(N \times N)}$ of pairwise Cosine similarities between documents. More specifically, the Laplacian matrix is computed as $L = D^{-1/2}AD^{-1/2}$, where D is a diagonal matrix. Each diagonal element contains the sum of the i -th row of similarities $D_{ii} = \sum_{j=1}^N A_{ij}$. The next step is the construction of a matrix $X^{(N \times k)} = \{x_i : i = 1, \dots, k\}$ whose columns correspond to the k

¹ <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-20/www/data/news20.tar.gz>.

² <http://www.daviddlewis.com/resources/testcollections/reuters21578/reuters21578.tar.gz>.

largest eigenvectors of L . The standard k -means algorithm is then used to cluster the rows of matrix X after being normalized to unit length in Euclidean space, where the i -th row is the vector representation of the i -th document in the new feature space.

Clustering evaluation was based on the supervised measure *Normalized Mutual Information (NMI)* and the F_1 -measure. We denote as n_i^{gt} the number of documents of class i , n_j the size of cluster j , n_{ij} the number of documents belonging to class i that are clustered in cluster j , C^{gt} the grouping based on ground truth labels of documents $c_1^{gt}, \dots, c_k^{gt}$ (true classes). Let us further denote $p(c_i^{gt}) = n_i^{gt}/N$ and $p(c_j) = n_j/N$ the probability of selecting arbitrarily a document from the data set and that belongs to class c_i^{gt} and cluster c_j , respectively, and $p(c_i^{gt}, c_j) = n_{ij}/N$ the joint of arbitrarily selecting a document from data set and that belongs to cluster c_j and is of class c_i^{gt} . Then, the [0,1]-Normalized MI measure is computed by dividing the Mutual Information by the maximum between the cluster and class entropy:

$$NMI(C^{gt}, C) = \left(\sum_{\substack{c_i^{gt} \in C^{gt} \\ c_j \in C}} p(c_i^{gt}, c_j) \log \frac{p(c_i^{gt}, c_j)}{p(c_i^{gt}) p(c_j)} \right) / \max\{H(C^{gt}), H(C)\}. \tag{16}$$

When C and C^{gt} are independent, the value of *NMI* equals to zero, while it equals to one if these partitions contain identical clusters.

The F_1 -measure is the harmonic mean of the *precision* and *recall* measures of the clustering solution:

$$F_1 = \frac{precision \cdot recall}{precision + recall}. \tag{17}$$

Higher values of F_1 in [0,1] indicate better clustering solutions.

Tables 2, 3, 5, and 6 present the results from the experiments conducted for each collection. Specifically, we compared the classic BOW representation, the GVSM, the proposed GTCV-VSM method (with $\lambda = 0.2$ in Eq. 11), that represents the documents as described in Eq. 14 and the CVM-VSM as proposed in [5], where document vectors are computed based on Eq. 6 with *idf* weights. More specifically, for each collection, each representation method was tested for 100 runs of spk-means (Tables 2, 3, 4) and spectral-c (Tables 5, 6, 7). To provide fair comparative results, for each document collection, all methods were initialized using the same random document seeds. The average of all runs (*avg*), the average of the worst 10% of the clustering solutions (*avg*_{10%}), and the best values are reported for each performance measure. The worst 10% concerns the 10% of the solutions with the lowest *Cohesion*, while the best clustering solution is that having the maximum *Cohesion* in the 100 runs (for spectral-c the sum of squared distances is considered for this purpose). The best result for each dataset is emphasized in each of the *avg* and *best* columns. Moreover, in Fig. 4, we present the average clustering performance of spk-means with respect to the value of λ parameter of Eq. 11 where, although not best for all cases, the value 0.2 we used seems to be a reasonable choice for all the data sets we have considered. Note that similar effect was observed for spectral-c method.

In order to illustrate the statistical significance of the obtained results, the well-known *t-test* was applied for each data set to determine the significance of the performance difference between our methods and the compared methods. We have considered the case where $\sigma = 10$ for the gaussian kernel for all data sets. Within a confidence interval of 95% and for the value of degrees of freedom equal to 198 (for two sets of 100 experiments each),

Table 2 *NMI* values of the clustering solution for VSM (BOW), GVSM, CVM-VSM and the proposed GTCV-VSM (for several values of σ) document representations using the spk-means algorithm

Method	σ	D_1			D_2			D_3			D_4			D_5		
		avg	best	avg10%	avg	best	avg10%	avg	best	avg10%	avg	best	avg10%	avg	best	avg10%
BOW	–	0.722	0.821	0.594	0.748	0.829	0.638	0.537	0.548	0.379	0.625	0.779	0.505	0.552	0.562	0.535
GTCV	1	0.749	0.854	0.601	0.767	0.845	0.638	0.544	0.564	0.372	0.667	0.793	0.515	0.570	0.578	0.561
	2	0.756	0.871	0.631	0.765	0.852	0.657	0.563	0.574	0.396	0.670	0.832	0.539	0.572	0.580	0.561
	5	0.773	0.881	0.687	0.777	0.864	0.662	0.577	0.602	0.400	0.688	0.851	0.539	0.589	0.633	0.578
	10	0.777	0.886	0.685	0.781	0.873	0.672	0.590	0.621	0.424	0.684	0.849	0.540	0.590	0.630	0.580
	30	0.761	0.879	0.659	0.776	0.863	0.653	0.579	0.590	0.369	0.683	0.842	0.518	0.576	0.612	0.568
	inf	0.760	0.862	0.631	0.772	0.862	0.639	0.574	0.586	0.366	0.681	0.840	0.521	0.576	0.610	0.566
GVSM	–	0.752	0.832	0.611	0.747	0.822	0.637	0.556	0.576	0.419	0.670	0.827	0.547	0.575	0.580	0.573
CVM	–	0.750	0.841	0.612	0.754	0.851	0.659	0.547	0.604	0.400	0.672	0.824	0.541	0.578	0.581	0.575

Table 3 F_1 -measure values of the spk-means clustering solution for the different representation methods

Method	σ	D_1			D_2			D_3			D_4			D_5		
		avg	best	avg10%	avg	best	avg10%	avg	best	avg10%	avg	best	avg10%	avg	best	avg10%
BOW	–	0.779	0.920	0.685	0.780	0.901	0.645	0.703	0.706	0.570	0.735	0.918	0.558	0.675	0.697	0.646
GTCV	1	0.806	0.940	0.688	0.790	0.921	0.650	0.709	0.713	0.576	0.755	0.920	0.561	0.691	0.695	0.677
	2	0.814	0.946	0.688	0.792	0.924	0.674	0.721	0.728	0.580	0.764	0.938	0.598	0.698	0.714	0.672
	5	0.828	0.953	0.722	0.817	0.929	0.665	0.736	0.737	0.597	0.773	0.948	0.611	0.712	0.751	0.681
	10	0.832	0.954	0.733	0.820	0.936	0.603	0.737	0.739	0.603	0.773	0.947	0.581	0.712	0.749	0.681
	30	0.814	0.950	0.747	0.794	0.929	0.657	0.725	0.727	0.576	0.766	0.944	0.579	0.698	0.746	0.666
	inf	0.813	0.942	0.689	0.792	0.926	0.651	0.722	0.728	0.576	0.765	0.944	0.581	0.698	0.744	0.666
GVSM	–	0.790	0.923	0.705	0.783	0.903	0.640	0.706	0.71	0.576	0.750	0.943	0.591	0.687	0.720	0.672
CVM	–	0.765	0.941	0.672	0.790	0.930	0.672	0.708	0.725	0.576	0.751	0.934	0.604	0.685	0.716	0.669

Table 4 The p and t values of the statistical significance t -test of the difference in k -means performance using GTCV-VSM ($\sigma = 10$) and the compared representation methods, with respect to the two evaluation measures

GTCV	D_1	D_2		D_3		D_4		D_5		
		p -val	t -val	p -val	t -val	p -val	t -val	p -val	t -val	
$(\sigma = 10)$ vs										
BOW _{NMI}	0.011×10^{-6}	5.98	0.075×10^{-3}	4.05	0.025×10^{-6}	5.81	0.080×10^{-8}	6.45	.0000	12.8
GVSM _{NMI}	0.00008	2.68	0.081×10^{-3}	4.02	0.050×10^{-3}	4.15	0.085	1.73	0.056×10^{-5}	5.17
CVM _{NMI}	0.0051	2.83	0.0010	3.33	0.052×10^{-4}	4.65	0.1659	1.39	0.077×10^{-3}	4.04
BOW _{F1}	0.020×10^{-5}	5.39	0.050×10^{-2}	3.54	0.046×10^{-2}	3.56	0.0010	3.32	0.0000	12.8
GVSM _{F1}	0.037×10^{-3}	4.22	0.00021	3.11	0.067×10^{-2}	3.45	0.0329	2.15	0.0000	9.06
CVM _{F1}	0.081×10^{-3}	4.02	0.06×10^{-8}	6.50	0.0027	3.04	0.0314	2.18	0.0000	9.31

Values of p smaller than the significance level of 0.05 (5%) indicate significant superiority of GTCV-VSM

Table 5 *NMI* values of the clustering solution for VSM (BOW), GVSM, CVM-VSM and the proposed GTCV-VSM (for several values of σ) document representations using the spectral clustering algorithm

Method	σ	D_1			D_2			D_3			D_4			D_5		
		<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>
BOW	–	0.753	0.761	0.750	0.781	0.788	0.737	0.569	0.585	0.555	0.718	0.780	0.631	0.558	0.559	0.506
GTCV	1	0.770	0.774	0.769	0.790	0.795	0.750	0.614	0.626	0.600	0.735	0.779	0.642	0.560	0.561	0.516
	2	0.781	0.785	0.760	0.790	0.794	0.757	0.625	0.632	0.601	0.752	0.789	0.649	0.562	0.564	0.523
	5	0.794	0.804	0.790	0.833	0.853	0.763	0.639	0.640	0.619	0.768	0.827	0.669	0.579	0.600	0.557
	10	0.807	0.814	0.801	0.833	0.853	0.761	0.645	0.648	0.620	0.758	0.819	0.661	0.581	0.589	0.558
	30	0.791	0.796	0.769	0.807	0.832	0.743	0.613	0.613	0.609	0.755	0.797	0.647	0.567	0.582	0.535
	inf	0.774	0.782	0.767	0.794	0.794	0.722	0.619	0.619	0.610	0.749	0.793	0.637	0.560	0.568	0.530
GVSM	–	0.756	0.770	0.702	0.794	0.830	0.747	0.593	0.595	0.586	0.722	0.780	0.637	0.548	0.554	0.513
CVM	–	0.761	0.768	0.751	0.801	0.823	0.760	0.605	0.606	0.590	0.728	0.794	0.642	0.557	0.566	0.519

Table 6 F_1 -measure values of the spectral clustering solution for the different representation methods

Method	σ	D_1			D_2			D_3			D_4			D_5		
		<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>	<i>avg</i>	<i>best</i>	<i>avg10%</i>
BOW	–	0.801	0.811	0.780	0.819	0.822	0.767	0.710	0.723	0.701	0.808	0.911	0.697	0.666	0.669	0.654
GTCV	1	0.811	0.819	0.809	0.822	0.832	0.772	0.729	0.741	0.728	0.834	0.915	0.722	0.694	0.703	0.663
	2	0.818	0.823	0.806	0.837	0.841	0.779	0.733	0.746	0.732	0.865	0.922	0.725	0.689	0.703	0.652
	5	0.837	0.840	0.818	0.887	0.927	0.792	0.744	0.756	0.737	0.870	0.930	0.740	0.716	0.727	0.647
	10	0.840	0.842	0.826	0.890	0.925	0.788	0.754	0.759	0.742	0.865	0.929	0.736	0.710	0.725	0.654
	30	0.823	0.826	0.809	0.856	0.886	0.769	0.726	0.735	0.725	0.864	0.925	0.705	0.704	0.701	0.642
	inf	0.814	0.817	0.806	0.826	0.832	0.734	0.728	0.735	0.729	0.859	0.922	0.703	0.692	0.686	0.653
GVSM	–	0.756	0.770	0.702	0.826	0.901	0.780	0.709	0.714	0.724	0.823	0.916	0.705	0.642	0.657	0.654
CVM	–	0.761	0.768	0.779	0.831	0.897	0.791	0.725	0.725	0.723	0.825	0.916	0.713	0.673	0.678	0.654

Table 7 The p and t values of the statistical significance t -test of the difference in spectral clustering performance using GTCV-VSM ($\sigma = 10$) and the compared representation methods, with respect to the two evaluation measures

GTCV ($\sigma = 10$) vs	D_1		D_2		D_3		D_4		D_5	
	<i>p-val</i>	<i>t-val</i>	<i>p-val</i>	<i>t-val</i>	<i>p-val</i>	<i>t-val</i>	<i>p-val</i>	<i>t-val</i>	<i>p-val</i>	<i>t-val</i>
BOW_{NMI}	0.0000	27.3	0.0000	13.8	0.0000	620	0.026×10^{-4}	4.85	0.0000	8.03
$GVSM_{NMI}$	0.0000	16.7	0.0000	7.51	0.0000	130	0.129×10^{-5}	4.99	0.0000	12.1
CVM_{NMI}	0.0000	19.3	0.150×10^{-8}	6.35	0.0000	138	0.316×10^{-3}	3.67	0.0000	8.83
BOW_{F_1}	0.0000	24.1	0.0000	11.4	0.0000	875	0.123×10^{-4}	4.48	0.0000	19.1
$GVSM_{F_1}$	0.0000	15.1	0.0000	7.53	0.0000	410	0.113×10^{-2}	3.31	0.0000	30.7
CVM_{F_1}	0.0000	18.7	0.0000	7.11	0.0000	268	0.115×10^{-3}	3.94	0.0000	14.1

Values of p smaller than the significance level of 0.05 (5%) indicate significant superiority of GTCV-VSM

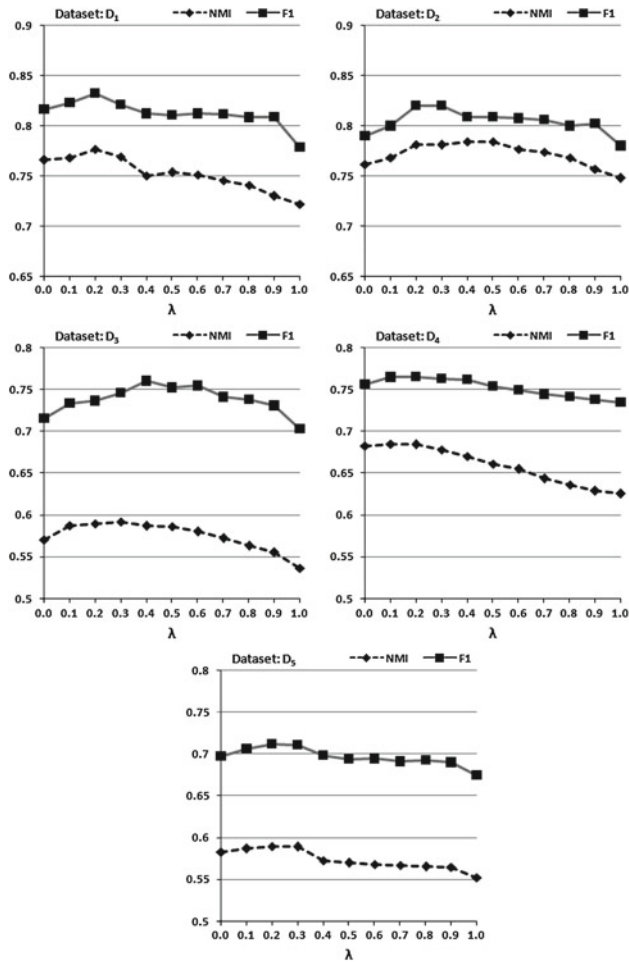


Fig. 4 The effect of varying the parameter λ on the spk-means clustering performance for each data set. Eq. 11 is used to determine the term self-weight α_v when computing the *l*tcv histograms

the critical value for t is $t_c = 1.972$ ($p_c = 5\%$ for p value). This means that if the computed $t \geq t_c$, then the null hypothesis is rejected ($p \geq 5\%$, respectively), i.e., our method is superior, otherwise the *null hypothesis* is accepted. As it can be observed from the results of the statistical tests for spk-means presented in Table 4, the performance superiority of GTCV-VSM is clearly significant in four out of five data sets with respect to all other methods. For data set D_4 , the tests indicate that GTCV-VSM, although still better than BOW, has less significant difference in performance compared with GVSM and CVM-VSM. Table 4 provides the respective t-test results for the spectral-c method where, also due to the lower standard deviation of the results using all document representation methods, the GTCV-VSM demonstrates significantly better results than the compared representations.

The experimental results indicate that our method outperforms the traditional BOW approach in all cases, even for small values of smoothing parameter σ (e.g., $\sigma = 1$ or 2). This substantiates our rationale that the clustering procedure is assisted by the proposed semantic

smoothing, which takes into account the local contextual information associated with a term occurrence. GTCV-VSM requires moderate values for the parameter σ to achieve better performance. The same is observed for the quality (in terms of NMI or F_1) of the best solution (i.e., the one with maximum *Cohesion*) found in the 100 runs, where moderate values of σ (i.e., $\sigma = 5$ or 10) result in better GTCV-VSM performance. Moreover, the clustering results for a wide range of values of the smoothing parameter σ indicate that the method is quite robust to the specification of this parameter. GTCV-VSM behaves similarly to BOW when a low value is set for σ , while when this value becomes very high, the discriminative information of the global term context vectors is reduced. This was demonstrated using spk-means and spectral clustering methods. Among them, the latter in all cases except from D_5 presented better average clustering solutions in terms of both evaluation measures NMI and F_1 , while interestingly, spk-means was superior in terms of the best clustering solutions in most cases (with the exception of D_3) despite operating in a feature space of a much larger size.

6 Conclusions

We have presented the global term context vector-VSM (GTCV-VSM) document representation, an extension to the vector space model that determines a proper feature space to project the typical VSM document vector representations. Our approach is entirely corpus-based and operates in the preprocessing phase in a sequence of four steps: (i) captures local contextual information associated with each term occurrence in the term sequences of documents; (ii) summarizes the local context vectors of each term into the respective global term context vectors; (iii) constructs the semantic matrix for a problem using the global term context vectors; and finally, (iv) projects documents using the semantic matrix. The proposed approach achieves semantic smoothing by reducing data sparsity, while retaining the original dimensionality. The derived representation maintains the initial interpretability since each dimension is associated with a single vocabulary term. In the experimental document clustering study, we compared the proposed representation with the typical VSM, the Generalized-VSM and CVM-VSM, using Cosine similarity. The statistical analysis of the obtained results indicates that GTCV-VSM assists well-known clustering algorithms, such as spherical k -means and spectral clustering, to achieve better clustering solutions compared with other representation methods.

Our plans for future work are to investigate the potential of combining the local and global contextual information associated with terms to explore ways of building compact concept vectors, to efficiently project the transformed document vectors in feature spaces of lower dimensionality, and to perform a systematic study for procedures that could efficiently compute α_v parameters (Eq. 13) for each vocabulary term, which could improve the global term context vectors. Finally, we aim at examining the proposed representation for document classification.

References

1. AlSumait L, Domeniconi C (2008) Text clustering with local semantic kernels. In: Berry M, Castellanos M (eds) Survey of text mining II. Springer, London, pp 219–232
2. Apté C, Damerau F, Weiss SM (1994) Towards language independent automated learning of text categorization models. In: SIGIR '94: proceedings of the 17th annual international ACM SIGIR conference on research and development in information retrieval. Springer, New York, pp 23–30

3. Beil F, Ester M, Xu X (2002) Frequent term-based text clustering. In: KDD '02: proceedings of the 8th ACM SIGKDD international conference on knowledge discovery and data mining. ACM, New York, pp 436–442. doi:[10.1145/775047.775110](https://doi.org/10.1145/775047.775110)
4. Beyer K, Goldstein J, Ramakrishnan R, Shaft U (1999) When is “nearest neighbor” meaningful? In: ICDT '99: proceedings of the 7th international conference on database theory. Springer, London, pp 217–235
5. Billhardt H, Borrajo D, Maojo V (2002) A context vector model for information retrieval. *J Am Soc Inf Sci Technol* 53(3):236–249. doi:[10.1002/asi.10032](https://doi.org/10.1002/asi.10032)
6. Chen C, Tseng F, Liang T (2010) An integration of fuzzy association rules and wordnet for document clustering. *Knowl Inf Syst* (available online). doi:[10.1007/s10115-010-0364-2](https://doi.org/10.1007/s10115-010-0364-2)
7. Deerwester S, Dumais S, Furnas G, Landauer T, Harshman R (1990) Indexing by latent semantic analysis. *J Am Soc Inf Sci* 41:391–407
8. Dhillon I, Modha D (2001) Concept decompositions for large sparse text data using clustering. *Mach Learn* 42(1):143–175. doi:[10.1023/A:1007612920971](https://doi.org/10.1023/A:1007612920971)
9. Farahat A, Kamel M (2010) Statistical semantics for enhancing document clustering. *Knowledge and Information Systems* (available online). doi:[10.1007/s10115-010-0367-z](https://doi.org/10.1007/s10115-010-0367-z)
10. Fung B, Wang K, Ester M (2003) Hierarchical document clustering using frequent itemsets. In: Proceedings of SIAM international conference on data mining
11. Ghosh J, Strehl A (2006) Similarity-based text clustering: a comparative study. In: Kogan J, Nicholas C, Teboulle M (eds) *Grouping multidimensional data*. Springer, Berlin, pp 73–97
12. Grauman K, Darrell T (2007) The pyramid match kernel: efficient learning with sets of features. *J Mach Learn Res* 8:725–760. doi:[10.1145/361219.361220](https://doi.org/10.1145/361219.361220)
13. Hotho A, Maedche E, Staab S (2001) Ontology-based text document clustering. *Knstliche Intell* 4:48–54
14. Hu X, Sun N, Zhang C, Chua T (2009) Exploiting internal and external semantics for the clustering of short texts using world knowledge. In: Proceeding of the 18th ACM conference on information and knowledge management. ACM, New York, CIKM '09, pp 919–928. doi:[10.1145/1645953.1646071](https://doi.org/10.1145/1645953.1646071)
15. Jing J, Zhou L, Ng M, Huang Z (2006) Ontology-based distance measure for text clustering. In: Proceedings SIAM SDM workshop on text mining
16. Karypis G, Han E (2000) Concept indexing: a fast dimensionality reduction algorithm with applications to document retrieval and categorization. In: Technical report TR-00-0016. University of Minnesota
17. Keikha M, Razavian N, Oroumchian F, Razi H (2008) Document representation and quality of text: an analysis. In: Berry M, Castellanos M (eds) *Survey of text mining II*. Springer, London, pp 219–232
18. Lebanon G, Mao Y, Dillon J (2007) The locally weighted bag of words framework for document representation. *J Mach Learn Res* 8:2405–2441
19. Lewis D (1992) An evaluation of phrasal and clustered representations on a text categorization task. In: SIGIR '92: Proceedings of the 15th annual international ACM SIGIR conference on research and development in information retrieval. ACM, New York, pp 37–50. doi:[10.1145/133160.133172](https://doi.org/10.1145/133160.133172)
20. Li Y, Chung S, Holt J (2008) Text document clustering based on frequent word meaning sequences. *Data Knowl Eng* 64(1):381–404. doi:[10.1016/j.datak.2007.08.001](https://doi.org/10.1016/j.datak.2007.08.001)
21. McQueen J (1967) Some methods for classification and analysis of multivariate observations. In: Proceedings of 5th Berkley symposium on mathematical statistics and probability. pp 281–297
22. Miller G, Beckwith R, Fellbaum C, Gross D, Miller K (1990) Wordnet: an on-line lexical database. *Int J Lexicogr* 3:235–244
23. Mladenic D (1998) Machine learning on non-homogeneous, distributed text data. PhD thesis, University of Ljubljana, Faculty of Computer and Information Science
24. Ng A, Jordan M, Weiss Y (2001) On spectral clustering: analysis and an algorithm. *Adv Neural Inf Process Syst* 14:849–864
25. Ni X, Quan X, Lu Z, Wenxin L, Hua B (2010) Short text clustering by finding core terms. *Knowl Inf Syst* 1–21. doi:[10.1007/s10115-010-0299-7](https://doi.org/10.1007/s10115-010-0299-7)
26. Porter M (1997) An algorithm for suffix stripping. In: Jones K, Willett P (eds) *Readings in information retrieval*. Morgan Kaufmann Publishers, San Francisco, pp 313–316
27. Pu W, Liu N, Yan S, Yan J, Xie K, Chen Z (2007) Local word bag model for text categorization. In: ICDM '07: proceedings of the 2007 7th IEEE international conference on data mining. IEEE Computer Society, Washington, pp 625–630. doi:[10.1109/ICDM.2007.69](https://doi.org/10.1109/ICDM.2007.69)
28. Salton G, Wong A, Yang C (1975) A vector space model for automatic indexing. *Commun ACM* 18(11):613–620. doi:[10.1145/361219.361220](https://doi.org/10.1145/361219.361220)
29. Wang P, Domeniconi C (2008) Building semantic kernels for text classification using wikipedia. In: KDD '08: proceeding of the 14th ACM SIGKDD international conference on knowledge discovery and data mining. ACM, New York, pp 713–721. doi:[10.1145/1401890.1401976](https://doi.org/10.1145/1401890.1401976)
30. Wikipedia (2004) Wikipedia, the free encyclopedia. <http://en.wikipedia.org/>

31. Wong S, Ziarko W, Wong P (1985) Generalized vector spaces model in information retrieval. In: SIGIR '85: proceedings of the 8th annual international ACM SIGIR conference on research and development in information retrieval. ACM, New York, pp 18–25. doi:[10.1145/253495.253506](https://doi.org/10.1145/253495.253506)

Author Biographies



Argyris Kalogeratos received the B.Sc. and M.Sc. degrees in Computer Science from the University of Ioannina, Ioannina, Greece, in 2006 and 2008, respectively. Currently, he is pursuing the Ph.D. degree in the Department of Computer Science, University of Ioannina. His research interests include machine learning, data clustering, text representation, and mining.



Aristidis Likas received the Diploma degree in electrical engineering and the Ph.D. degree in electrical and computer engineering from the National Technical University of Athens, Greece, in 1990 and 1994, respectively. Since 1996, he has been with the Department of Computer Science, University of Ioannina, Greece, where he is currently an Associate Professor. His research interests include machine learning, data mining, multimedia content analysis and bioinformatics.