

Visualization of Text Document Corpus

Blaž Fortuna, Marko Grobelnik and Dunja Mladenić
 Jozef Stefan Institute
 Jamova 39, 1000 Ljubljana, Slovenia
 E-mail: {blaz.fortuna, marko.grobelnik, dunja.mladenic}@ijs.si

Keywords: text visualization, latent semantic indexing, mutlidimensional scaling

Received: June 2, 2005

Visualization is commonly used in data analysis to help the user in getting an initial idea about the raw data as well as visual representation of the regularities obtained in the analysis. In similar way, when we talk about automated text processing and the data consists of text documents, visualization of text document corpus can be very useful. From the automated text processing point of view, natural language is very redundant in the sense that many different words share a common or similar meaning. For computer this can be hard to understand without some background knowledge. We describe an approach to visualization of text document collection based on methods from linear algebra. We apply Latent Semantic Indexing (LSI) as a technique that helps in extracting some of the background knowledge from corpus of text documents. This can be also viewed as extraction of hidden semantic concepts from text documents. In this way visualization can be very helpful in data analysis, for instance, for finding main topics that appear in larger sets of documents. Extraction of main concepts from documents using techniques such as LSI, can make the results of visualizations more useful. For example, given a set of descriptions of European Research projects (6FP) one can find main areas that these projects cover including semantic web, e-learning, security, etc. In this paper we describe a method for visualization of document corpus based on LSI, the system implementing it and give results of using the system on several datasets.

Povzetek: Predstavljena je vizualizacija korpusa besedil.

1 Introduction

Automated text processing is commonly used when dealing with text data written in a natural language. However, when processing the data using computers, we should be aware of the fact that many words having different form share a common or similar meaning. For a computer this can be difficult to handle without some additional information -- background knowledge. Latent Semantic Indexing (LSI) is a technique for extracting this background knowledge from text documents. It employs a technique from linear algebra called Singular Value Decomposition (SVD) and the bag-of-words representation of text documents for extracting words with similar meanings. This can also be viewed as the extraction of hidden semantic concepts from text documents.

Visualization of a document corpus is a very useful tool for finding the main topics that the documents from this corpus talk about. Different methods were proposed for visualizing a large document collection using different underlying methods. For instance, visualization of large document collection based on document clustering [3], or visualization of news collection based on visualizing relationships between named entities extracted from the text [4]. Another example used in our work is visualization of European research space [5].

Given a set of descriptions of European research projects in IT (6th Framework IST), using document visualization one can find main areas that these projects cover, such as *semantic web, e-learning, security*, etc.

In automated text processing document are usually represented using the bag-of-words document representation, where each word from the document vocabulary stands for one dimension of the multidimensional space of documents. Consequently, in automated text processing we are dealing with very high dimensionality of up to hundreds of thousands dimensions. Dimensionality reduction [6] is important for different aspects of automated text processing including document visualization.

We propose to use dimensionality reduction for document visualization by first extracting main concepts from documents using LSI and than using this information to position documents on a two dimensional plane via multidimensional scaling [1]. The final output is graphical presentation of a document set that can be plotted on a computer screen. The proposed approach is implemented as a part of *Text Garden* software tools for text mining [7]¹ in a component providing different kinds of document corpus visualization based on LSI and multidimensional scaling.

¹ <http://www.textmining.net/>

This paper is organized as follows. Section 2 provides a short description of LSI and multidimensional scaling, while its application to document visualization is given in Section 3. Description of the developed system implementing the method is given in Section 4. Section 5 provides conclusions and discussion.

2 Building Blocks

First step of our approach to visualization of a document corpus is mapping all the documents into two dimensional space so we can plot them on a computer screen. Ideally they would be positioned in such a way that the distance between two documents would correspond to the content similarity between them.

We obtain this mapping by sending the document corpora through the pipeline for reducing dimensionality, consisting from building blocks presented in this Section. The whole pipeline will be outlined in the Section 3.

2.1 Representation of Text Documents

The first step in our approach is to represent text documents as vectors. We use the standard Bag-of-Words (BOW) representation together with TFIDF weighting [9]. In the BOW representation there is a dimension for each word; a document is encoded as a feature vector with word frequencies as elements. Elements of vectors are weighted, in our case using the standard TFIDF weights as follows. The i -th element of the vector containing frequency of the i -th word is multiplied with $IDF_i = \log(N/df_i)$, where N is total number of documents and df_i is document frequency of the i -th word (the number of documents from the whole corpus in which the i -th word appears).

2.2 Latent Semantic Indexing

A well known and used approach for extracting latent semantics (or topics) from text documents is Latent Semantic Indexing [2]. In this approach we first construct term-document matrix A from a given corpus of text documents. This is a matrix with vectors of documents from a given corpus as columns. The term-document matrix A is then decomposed using singular value decomposition, so that $A = USV^T$; here matrices U and V are orthogonal and S is a diagonal matrix with ordered singular values on the diagonal. Columns of matrix U form an orthogonal basis of a subspace in the bag-of-words space where vectors with higher singular values carry more information -- this follows from the basic theorem about SVD, which tells that by setting all but the largest k singular values to 0 we get the best approximation for matrix A with matrix of rank k). Vectors that form the basis can be also viewed as concepts and the space spanned by these vectors is called the *Semantic Space*.

Each concept is a vector in the bag-of-words space, so the elements of this vector are weights assigned to the words coming from our documents. The words with the

highest positive or negative values form a set of words that are found most suitable to describe the corresponding concept.

A related approach (not used here) that also aims at extracting latent semantics from text documents is Probabilistic Latent Semantic Analysis (PLSA) introduced in [8]. Compared to standard Latent Semantic Analysis which comes from linear algebra and performs a Singular Value Decomposition of co-occurrence tables, this method is based on a mixture decomposition derived from a latent class model. This method assigns each word a probability to be in a concept, where the number of concepts is predefined.

2.3 Dimensionality Reduction

We are using a sequential combination of linear subspace methods and multidimensional scaling for reducing document space dimensionality. Both methods can be independently applied to any data set that is represented as a set of vectors in some higher dimensional space. Our goal is to lower the number of dimensions to two so that the whole corpus of documents can be shown on a computer screen.

Linear subspace methods [10], like Principal Component Analysis (PCA) or Latent Semantic Indexing, focus on finding direction in original vector space, so they capture the most variance (as is the case for PCA) or are the best approximation for original document-term matrix (as is the case for LSI). By projecting data (text documents) only on the first two directions we can get the points that live in the two dimensional space. The problem with linear subspace methods is that only the information from the first two directions is preserved. In case of LSI it would mean that all documents are described using only the two main concepts.

Multidimensional scaling [1] enables dimensionality reduction by mapping original multidimensional vectors onto two dimensions. Here the points representing documents are positioned into two dimensions so they minimize some energy function. The basic and most common form of this function is

$$E = \sum_{i \neq j} \delta_{ij} - d(x_i, x_j)^2,$$

where x_i are two dimensional points and δ_{ij} represents the similarity between two vectors (in our case documents i and j). An intuitive description of this optimization problem is: the better the distances between points on the plane approximate real similarity between documents, the lower the value of the energy function. Function E is nonnegative and equals zero only when distances between points match exactly with similarity between documents.

- [9] Salton, G. Developments in Automatic Text Retrieval, Science, Vol 253, pages 974-979, 1991.
- [10] Shawe-Taylor, J., Cristianini, N. Kernel Methods for Pattern Analysis, Cambridge University Press, 2004, 143-150
- [11] Vallbe, J.J., Marti, M.A., Fortuna, B., Jakulin, A., Mladenec, D., Casanovas, P. Stemming and lemmatisation, Springer Lecture Notes, (to appear), 2006.