

# Various Semantic based Models and Text Analysis- A Review

Swati Gautam, Hansa Acharya, Prof. Anurag Jain

**Abstract**— In this paper is to here, a lot of topics cover up in a document are more frequently than not related to the context of the document, evaluate topical ideas within circumstance can potentially make known many interesting theme patterns in classify to progression of text precisely and efficiently. In the text analysis, typical algorithms are not good enough for the reason that requisites frequently were assumed self-sufficient, which was be acquainted with as an extraction model. But, the extraction model suggests a relatively make poor illustration of the data for the reason that it compensate no attention to any correlations between the conditions.

**Index Terms**— Information Retrieval, Natural Language Processing.

## 1 INTRODUCTION

The indispensable confirmation of a textual matching model has to be recognized in unambiguous functional domains that make available methodologies and metrics for assessment. The rising quantity of textual data obtainable in electronic appearance is an main motivation for the search of well-organized techniques in the wide-ranging area of textual data looking are at in particular, Information Retrieval (IR). The most important purpose of IR is to proficiently recognize pertinent documents in a database, assuring an information could do with communicated by a user in a type of a query. This task to be converted into extra complicated as the size of the investigated databases enlarges and come within reach of aspiring at shrinking the size of the look for space by arrangement the document gatherings are currently also think about. In the domains of IR and textual database arrangement the design of well-organized textual similarities is a essential issue. In IR, the objective can be analyzed as the search, in a particular semantic space, of the documents for the most part similar to the query. And look for be capable of therefore be approved out through the working out of textual comparison between the query and each of the documents in the database. Document collected works arrangement can also be accomplished by clustering the documents that become visible close according to a well preferred, semantically argument, textual similarity.

The usual models make use of in Information Retrieval are found on a vector representation of the documents join together with a correspondence compute operation on the fundamental vector space. This model introduces a representation of the documents that put together extra "semantic" information by means of a sharing illustration of the semantics of the words. A text document is repeatedly join together with a variety of category of context information, such as the instance and position at which the document was generated, the author's who wrote the text, and its distributor. The substances of text documents with the comparable or similar context are frequently show a relationship in a number

of way. For ex:-news articles written in the period of a number of most important incident all have a tendency to be influenced by the occasion in some way, and papers written by the same researcher have a tendency to allocate comparable issues. In order to make known important content patterns in such contextualized text data, it is indispensable to be think about context information when investigating the areas covered in such data. Certainly, there have been more than a few latest studies in this direction. For example, the time stamps of text documents have been think about in some current work on sequential text mining [11, 18, 16,6]. Furthermore, author-topic analysis is considered in [19], and cross compilation proportional text mining is studied in [20]. All these studies consider some kinds of context information, i.e., time, authorship, and sub collected works. On the other hand, existing techniques are more often than not adjusted for some specific tasks, and are not appropriate to think about other category of contexts. For example, one cannot straightforwardly use the temporal text mining techniques to model the activity of authors. This is to be a sign of a serious drawback of existing contextual analysis of premises: every time when a new arrangement of context information is to be think about, people have to search for solutions in an ad hoc way.

Consequently, it is extremely attractive to commence a general text mining difficulty, contextual text mining, which is nonrepresentational from a family unit of text mining tasks with a choice of types of contextual analysis. It is attractive to originate a model that is very much common to accomplish the common tasks of these specific contextual text mining problems, and straightforward to be functional to each of them with suitable expected them.

In this work, we define the widespread difficulty of Contextual Text Mining (CtxTM) and its widespread tasks, which is nonrepresentational from a family unit of a particular text mining difficulty. We expand the probabilistic latent semantic analysis (PLSA) model to have as a feature context

information, and increase a contextual probabilistic latent semantic analysis (CPLSA) model to make possible contextual text mining in a wide-ranging way. By fitting the model to the text data to mine, we know how to (1) determine the comprehensive most important arguments from the gathering of documents; (2) evaluate the substance dissimilarity of the premises in some specified analysis of context; and (3) evaluate the exposure of themes join together with any particular circumstance.

These tasks are wide-ranging and can be simply applied to different specific contextual text mining problems. In this paper, they shows that a lot of stay alive contextual argued analysis problems can be defined as unique cases of CtxTM, and can be solved with make conformed adaptations of the combination model here they proposed, equivalent to the context information and the mining tasks it involves. Although it may not be the only promising model for contextual text mining, the model is reasonably flexible to adapt different assumptions.

## 2 PRIOR LITERATURE WORK

The image-based characterization of a objective word is a vector description co-event calculates with all illustration words transversely all the descriptions that enclose the objective word as a tag. The text-based and image-based vectors are concatenated and return to normal, consequence in the multimodal vectors that occupy the MDSM.

In Bruni et al. [2], here they estimate our own MDSM on the assignments of predicting semantic correspondence ruling, perception classification and imprison semantic neighbours of unusual classes. A careful explanation of our effects is that adding together image-based features is at least not destructive, when measure up to adding additional text-based features, and probably valuable. Prominently, in all experiments here they come across that image-based features show the way to concentrating qualitative discrepancy in concert. The MDSM is enhanced at imprisoning resemblance involving existing conceptions and focuses on their more image able properties (such as colour), while a equivalent text-based DSM is additional mechanism in the direction of nonrepresentational theory and belongings.

In paper[21], the correlated topic model (CTM) unambiguously models the correspondence between the concealed subject matters in the gathering and make possible the issue graphs creation and document browsers which consent to a user to find the way of collected works in a area conducted approach.

The correlated topic model manufactures on the earlier latent dirichlet allocation (LDA) model of [21] which is an example of a wide-ranging family unit of miscellaneous attachment models for decomposing data into multiple latent

components. LDA for the most part thinks that the words of each document take place from a combination of areas where each area is a multinomial in excess of a predetermined word vocabulary. The topics are collective by documents in the collected works but the topic shares change stochastically transversely documents as they are arbitrarily drawn from a Dirichlet distribution. A modern work in [21] has used LDA as a element is more complicated topic models. They not succeed to straightforwardly model the relationship between topics in the document.

In this paper[7], here they revise the CTM problem and propose a generative probabilistic mixture model for CTM. This model at the same time performs cross-collection clustering and within collected works clustering, and can be helpful to an individual set of comparable text collected works. The combinational model is foundation on constituent multinomial distribution models, each distinguish a different theme. The common arguments and collection specific themes are unambiguously modeled. The proposed model can be anticipated proficiently by means of the Expectation Maximization (EM) algorithm. We estimate [7]the model on two different text data sets i.e., a news article data set and a laptop analysis data set, and matched up to it with a baseline clustering technique also supported on a mixture model.

Distributional semantic models (DSMs)[13] ;fairly accurate the significance of words with vectors that stay behind follow of the samples, of co-incident of the words in a corpus, less than the assumption that semantically correlated words be supposed to come about in comparable contexts the distributional hypothesis [24].An even though their remarkable experimental accomplishment, DSMs are not exclusively reasonable as mental models of how we human beings get your hands on and make use of semantic acquaintance, in view of the fact that it is noticeable that we can rely not only on linguistic context, but also on our prosperous perceptual understanding[4]. In our existing research, we take up a extensive observation of distributional semantics. We put forward that word importance can be fundamentally captured by vectors summarizing co-occurrence patterns, but examine that co-occurrence necessitates not be inadequate to linguistic contexts. In exacting, our multimodal distributional semantic model (MDSM) take advantage of both co-occurrence with words on or after a average text corpus and co-occurrence with illustration characteristics take out using computer vision techniques from assortments of labeled images.

While present is a large literature by means of combined text and image-extracted characteristics to get better the categorization or classification of images [25], [23] we are in a intelligence shadow the opposite goal, i.e. take advantage of images to get better our rough calculation to word meaning. We are not more than responsive of an additional team practice

this goal, namely Feng and Lapata [14]. On the other hand, their contemporary model necessitates the taking out of a single distributional model from the same mixed media corpus, and text and visual words are characterized in expressions of the same collective latent dimensions. This has two significant drawbacks: First, the textual model be required to be removed from the same corpus images are taken from, and the text context extraction methods must be well-matched by means of the on the whole multimodal approach. Thus, image features cannot be additional to a state-of-the-art text-based DSM, to evaluate whether visual information is facilitating still when purely textual features are by now very good. Second, it is hard to evaluate the separate effect of image-based features on the overall performance in a joint model with concealed dimensions that mix textual and visual information. While standard SIFT-based visual words are good at characterizing parts of objects, we are currently augmenting our image-based vectors to include a richer set of visual features that capture properties such as texture, edges and color[15].

Consecutively to prevail over these concerns, they proposed a to some extent uncomplicated approach, in which the text-and image-based models are separately build from different sources, and then concatenated. We acquire the text-based distributional vector correspond to a word from a state-of-the-art DSM [12]and concatenate it with a vector that maintains follow of co-occurrence of the equivalent word with visual characteristics, take out from all the images in a make image collection that surround the word as a label (we use the ESP Game collection[17]. The visual features are, for the time being, SIFT-based visual words, a typical way, in computer vision, to represent images in expressions of counts of the discrete local interest point categories they surrounded [22].

### 3 FORMAL LANGUAGE PROCESSING

**NLP Processing and Ambiguity:** Natural Language Processing" (NLP) [9] as a control has been rising for many years. It was outward appearance in 1960 as a sub-field of Artificial Intelligence and Linguistics, with the endeavor of studying difficulty in the automatic generation and considerate of natural language. NLP is a very subsist and developing field of linguistics, with its many confronts still to defeat due to natural language's ambiguity. We have compensated special consideration to the distinctions between statistical and linguistic methods in natural language processing. Even the scientific communities that maintain each draw near are more often than not at odds, and NLP is frequently useful by using a arrangement of techniques from both approaches. Our experience in this area has made us wrap up that it is not probable to maintain that one approach is better than the other; this even includes the use of a mixed approach. Natural language, implicit as a tool that people utilize to express themselves, has explicit belongings that reduce the effectiveness of textual information recovery schemes. These

belongings are linguistic dissimilarity and uncertainty. By linguistic difference we denote the possibility of using different words or expressions to exchange a few words the equivalent scheme. Linguistic ambiguity is at what time a word or phrase allows for additional than one interpretation. Ambiguity, on the other hand, involve document noise, or the enclosure of non-meaningful documents, in view of the fact that texts were recovered that make use of the same term but with a different meaning. These distinctiveness put together automated language processing significantly complicated.

**Words:** Words is a part of language that bring sense of that thing. And it is the smallest component or unit make use of in content analysis. It is bring into play to generating result in a frequent allocation of a particular words or term.

**Parsing:** Parsing is a fundamental process[8] in any natural language processing pipeline, from the time when obtaining the syntactic arrangement of sentences make available us with information that can be make use of to take out meaning from them: essentials communicate to elements of meaning, and dependency relations illustrate the ways in which they work together, such as who carry out the action expressed in a sentence or which objective is receiving the action. Consequently, we can find parsers applied to many convenient difficulty in natural language processing where a few degree of semantic analysis is essential or well-situated, such as information extraction, information retrieval, machine conversion, textual demanding, or question answering.

### 4 DOCUMENT PROCESSING

**Language in the Electronic Age:** Content analysis is a creation of the electronic age[10]. Despite the fact that content analysis was frequently carry out in the 1940s, it became a more realistic and repeatedly used research techniques since the mid-1950's, as researchers initiated to center of attention on conceptions to a certain extent than simply words, and on semantic relationships relatively than just presence (de Sola Pool, 1959). Content analysis is a research tool centered on the genuine content and internal characteristics of media. It is used to establish the occurrence of certain words, phrases, characters, theory, topics, or sentences within texts or sets of texts and to compute this occurrence in an intent manner. Texts can be distinct generally as books, book chapters, interviews, discussions, newspaper headlines ,essays, and piece of writings, past manuscripts, verbal communications, discussions, marketing, show business, informal chat, or in authentic whatever thing happening of forthcoming language. To act upon a content investigation on a text, the text is set of laws, or not working down, into well-located types on a variety of levels word, word sense, phrase, sentence, or topic and then observing using one of content analysis basic methods: theoretical analysis or relational analysis. The consequences are

then make use of to make implications about the messages within the audience, the text's the writer's and yet the way of life and time of which these are a component.

**Information Warfare:** Definition: "the exploit and management of information in search of a spirited development in excess of an contender." Email spam, link spam, etc. Entire websites are makes up with counterfeit contented. Spammers by means of social networks to personalize attacks. BBC reports expectation in information on the web is individual damaged by the huge numbers of public paid by companies to post comments" (Dec. 2011). Public visualize news organizations have a susceptibility to sustain one side, and\are frequently operated by trustworthy people and organizations" declares by Pew Research Center (Sept. 2011).

**Why Analyses Documents:** Due to the reality that it can be functional to observe any part of a set of writing or incidence of evidence communication content analysis is make use of in large number of fields, ranging from advertising and media studies, to literature and expression, ethnography and intellectual learning, gender and age concerns, sociology and political science, psychology and cognitive science, as well as other fields of investigation[10]. Furthermore, content analysis replicates a close relationship with socio and psycholinguistics, and is participating an vital role in the development of artificial intelligence. The following list (personalized from Berelson, 1952) recommends more potential for the make use of content analysis:

- Expose international diversity in communication content.
- Identify the reality of misinformation.
- Identify the meanings, center of attention or communication developments of an individual, group or institution.
- Illustrate approach and performance reply to communications.
- Verify emotional or expressive circumstances of persons or groups.

## 5 DOCUMENT ANALYSIS

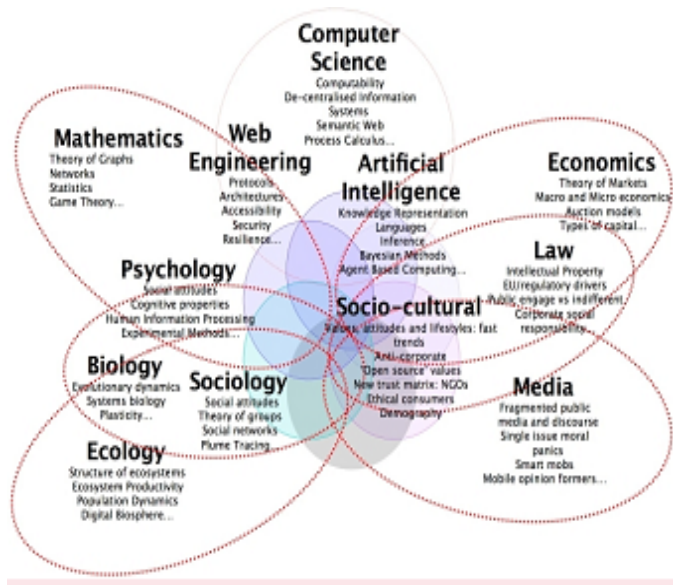
**Representation:** Document representation has a large impact on the performance of document retrieval and clustering algorithms. In view of the fact that contented words can become together into semantic classes nearby has been a considerable interest in low-dimensional term and document representations. Since any classifier is unable to understand a document in its unrefined format, a document has to be exchanged into a normal representation. Extensive work is carried out to propose various text representation techniques and text classification methods in the literature. But, it is

essential for researchers/practitioners to have a complete knowledge on all existing representation schemes and classifiers in order to select an appropriate representation scheme and classifier which best suits their purpose for an application.

**Resources:** In Text Mining learning, a document is normally bring into played as the necessary component of analysis. A *document* is a sequence of words and punctuation, subsequent the grammatical laws of the verbal communication, containing any relevant segment of text and can be of any length. It can be the document, an article, book, web page, emails, etc, depending on the type of analysis being performed and depending upon the goals of the researcher. In some cases, a document may surround no more than a chapter, a single paragraph, or even a single sentence. The fundamental unit of text is a *word*. A *term* is usually a word, but it can also be a word-pair or phrase. In this thesis, we will use *term* and *word* interchangeably. *Words* are consist of characters, and are the fundamental components from which meaning is constructed. By come together a word by means of grammatical structure, a sentence is prepared. *Sentences* are the fundamental element of accomplishment in text, containing information about the action of some subject. *Paragraphs* are the fundamental unit of composition and contain a related series of ideas or actions. As the length of text increases, additional structural forms become relevant, often including sections, chapters, whole documents, and to end with, a corpus of documents. A *corpus* is a collected works of documents. And, a *lexicon* is the set of all distinctive words in the corpus.

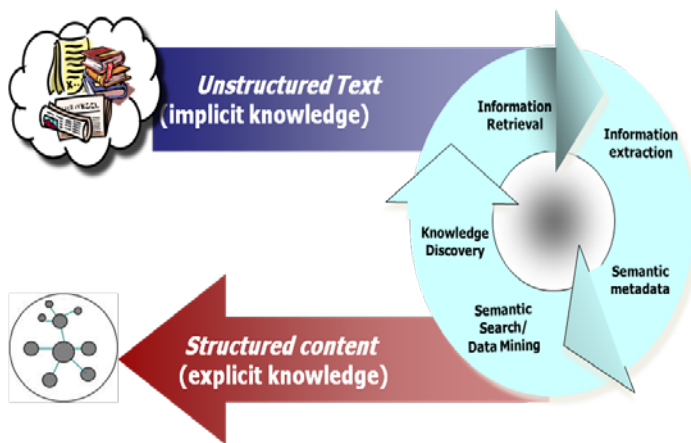
**Other Areas:** Additionally, like any new research method, content analysis do the accepted things to three fundamental principles of methodical technique that are as follows:

- Impartiality: Which represents that the analysis is followed on the foundation of unambiguous rules, which make possible unusual researchers to get hold of the equivalent results from the identical documents or messages.
- Organized: The inclusion or exclusion of substance is done according to a number of constantly useful rules where by the opportunity of together with only materials which sustain the researcher's ideas – is do away with.
- Simplified: The consequences get hold of by the researcher can be applied to other similar situations.



## 6 TEXT MINING & ITS TECHNIQUES

Text mining refers to the automated extraction of knowledge and information from text by means of revealing relationships and patterns present, but not obvious, in a document collection[3]. Text mining covers a broad field of tasks including text clustering, text categorization, information extraction, document summarization, sentiment analysis, named entity recognition, question answering and is an interdisciplinary field based on artificial intelligence, machine learning and statistics, natural language processing, information retrieval, computational linguistics, data mining. [1]



Text Mining is an interdisciplinary field that utilizes techniques from the general field of Data Mining and furthermore, joins methods from a variety of extra regions such as Information Retrieval, Computational Linguistics, Summarization, Information Extraction, Categorization, Clustering, Topic Tracking and Concept Linkage[5]. In the following sections, we will discuss each of these technologies and the role that they play in Text Mining.

**Information Extraction:** Information extraction (IE) is a process of automatically extracting structured information from unstructured and/or semi-structured machine-readable documents, processing human language texts by means of NLP[5]. The final output of the extraction process is some type of database obtained by looking for predefined sequences in text, a process called pattern matching. Tasks performed by IE systems include:

- *Expression analysis*, which identifies the expressions become visible in a document. This is particularly valuable for documents that restrain many difficult multi-word expressions, such as scientific research papers.
- *Named-unit identification*, which identifies the names appearing in a document, such as names of people or organizations. a number of structures are also able to be familiar with dates and expressions of time, capacity and percentages, combined elements, and so on.
- *Piece of information extraction*, which identifies and extracts complex facts from documents. Such piece of information could be links between entities or events.

b) **Information Retrieval:** Retrieval of text-based information also termed Information Retrieval (IR) has become a topic of great interest with the advent of text search engines on the Internet[5]. Text is considered to be composed of two fundamental units, namely the document (book, paragraphs, chapters, subdivisions, journal paper, Web pages, computer source code, and so onwards) and the term (word, word-pair, and phrase contained by a text). Conventionally in IR, text queries and manuscripts both are characterized in a combined approach, as locates of terms, to calculate the spaces between queries and documents accordingly providing a structure surrounded by directly implement simple text retrieval algorithms.

c) **Computational Linguistics/ Natural Language Processing:** Natural Language Processing [5] is a theoretically motivated range of computational techniques for analyzing and representing unsurprisingly happening texts at one or more levels of linguistic analysis for the intention of accomplishing human-like language processing for a collection of assignments or applications. The objective of Natural Language Processing (NLP) is to intend and make a computer system that will evaluate, identify with, and make natural human-languages. Purpose of NLP contain machine conversion of one human-language text to a different; creation of human-language text such as narrative, instruction manuals, and wide-ranging explanations; interfacing to other systems such as databases and robotic systems therefore make possible the use of human-language type commands and queries; and considerate human-language text to make available a review or to draw conclusions. NLP method gives the following tasks:

- Parse a sentence to find out its syntax.
- Find out the semantic importance of a sentence.
- Investigate the text context to agreed on its true meaning for measure up to it with other text.

**d) Categorization:** Categorization is the process of be acquainted with make a distinction and considerate the inspirations and things to group them into categories, for specific purpose[5]. In an ideal world, a type of clarifies a connection between the subjects and objects of knowledge. Categorization is essential in language, assumption, forecast, decision making and in all types of environmental relations. There are a lot of classification theories and techniques. In a broader past analysis, conversely three wide-ranging move towards to classification may be acknowledged as:

- *Conventional categorization* -According to the conventional analysis, a class be supposed to be noticeably distinct, mutually exclusive and cooperatively extensive, be in the right place to one, and only one, of the proposed type.
- *Theoretical clustering* - It is a present discrepancy of the traditional approach in which classes (clusters or entities) are created by first prepare their theoretical explanations and then categorizing the entities according to these explanations. Theoretical clustering is intimately communicated to fuzzy set theory, in which objects may fit in to one or more groups, in untrustworthy quantity of condition.
- *Sample theory* -Categorization can also be analysis as the procedure of grouping things based on prototypes. Categorization based on prototypes is the foundation for human development, and relies on be trained about the world passing through incarnation.

**e) Concept Linkage:** Concept linkage identifies related documents based on commonly shared concepts and between them[5]. The primary goal of concept linkage is to provide browsing for information rather than searching for it as in IR. For example, a Text Mining software solution may easily identify a link between topics X and Y, and Y and Z. Concept linkage is a valuable concept in Text Mining which could also detect a potential link between X and Z, rather than a human researcher has not appear transversely because of the large volume of information he or she would have to variety from beginning to end to create the relation. Conception relationship is valuable to classify relationships between diseases and handlings. In the in close proximity to expectations, Text Mining tools with concept linkage qualifications will be favorable in the biomedical field facilitating researchers to determine new treatments by associating treatments that have been used in interrelated fields.

## 7 CONCLUSION

In this paper, we define and revise a work on literature text mining trouble pass on to as proportional text mining. It has to do with find out any concealed frequent premises transversely a set of equivalent collected works on text as well as summarizing the relationship and differences of these collections beside each frequent premise. Existing come within reach of to sustain users in the search process do not deliberate

on the level of semantic matching required for searching concepts data on Web. This paper presents a concise preface to the different text representation proposals and classifiers utilized in the field of text mining. The obtainable techniques are measure up to and dissimilarity found on a variety of constraints to be precise criteria used for classification, algorithms take up and classification time complication. Since the above argument it is implicit that no single representation method and classifier can be suggested as a common representation for any application.

## REFERENCES

- [1] Vidhya & Aghila, "application of text mining and its relation to other fields and techniques", 2010.
- [2] Elia Bruni, Giang Binh Tran, and Marco Baroni. Distributional semantics from text and images. In Proceedings of the EMNLP GEMS Workshop, Edinburgh, 2011. In press.
- [3] SOPHIA ANANIADOU, School of Computer Science Director, National Centre for Text Mining. [www.nactem.ac.uk](http://www.nactem.ac.uk)
- [4] Max Louwerse. Symbol interdependency in symbolic and embodied cognition. Topics in Cognitive Science,3:273-302, 2011.
- [5] Winter School on "Data Mining Techniques and Tools for Knowledge Discovery in Agricultural Datasets" pg. no: 357-363.
- [6] C. C. Chen, M. C. Chen, and M.-S. Chen. Liped: Hmm-based life profiles for adaptive event detection. In Proceeding of KDD '05, pages 556-561, 2005.
- [7] ChengXiang Zhai, Atulya Velivelli, Bei Yu, "A Cross-Collection Mixture Model for Comparative Text Mining" KDD'04, Seattle, Washington, USA. ACM 1-58113-888,2004.
- [8] Carlos Gómez-Rodríguez," Mathematics, Computing, Language, and Life: Frontiers in Mathematical Linguistics and Language Theory: Volume 1 Parsing Schemata for Practical Text Analysis " <http://www.worldscibooks.com/compsci/p714.html>
- [9] <http://www.upf.edu/hipertextnet/en/numero-5/pln.html>
- [10] <https://www.ischool.utexas.edu/~palmquis/courses/content.html>
- [11] J. Kleinberg. Bursty and hierarchical structure in streams. In Proceedings of KDD '02, pages 91-101.
- [12] Marco Baroni and Alessandro Lenci. Distributional Memory: A general framework for corpus-based semantics. Computational Linguistics, 36(4):673-721, 2010.
- [13] Peter Turney and Patrick Pantel. From frequency to meaning: Vector space models of semantics. Journal of Artificial Intelligence Research, 37:141-188, 2010.
- [14] Yansong Feng and Mirella Lapata. Visual information in semantic representation. In Proceedings of HLT-NAACL, pages 91{99, Los Angeles, CA, 2010.
- [15] Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. In Proceedings of CVPR, pages 1778{1785, Miami Beach, FL, 2009.
- [16] Q. Mei and C. Zhai. Discovering evolutionary theme patterns from text: an exploration of temporal text mining. In Proceeding of KDD'05, pages 198-207,2005.
- [17] Luis von Ahn and Laura Dabbish. Labeling images with a computer game. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, pages 319{326, Vienna, Austria, 2004.
- [18] J. Perkiö, W. Buntine, and S. Perttu. Exploring independent trends in a topic-based search engine. In Proceedings of WI '04, pages 664-668, 2004.

- [19] M.Steyvers, P.Smyth, M.Rosen-Zvi, and T. Griffiths. Probabilistic author-topic models for information discovery. In Proceedings of KDD'04, pages 306–315, 2004.
- [20] C. Zhai, A. Velivelli, and B. Yu. A cross-collection mixture model for comparative text mining. In Proceedings of KDD'04, pages 743–748, 2004.
- [21] Blei, D. and Laerty, D.J. 2008. A correlated topic model of science. *Ann. Appl. Stat.* 1: 17-35.
- [22] Jun Yang, Yu-Gang Jiang, Alexander G. Hauptmann, and Chong-Wah Ngo. Evaluating bag-of-visual-words representations in scene classification. In *Multimedia Information Retrieval*, pages 197-206, 2007.
- [23] Jonathon Hare, Sina Samangooei, Paul Lewis, and Mark Nixon. Semantic spaces revisited: Investigating the performance of auto-annotation and semantic retrieval using semantic spaces. In *Proceedings of CIVR*, pages 359-368, Niagara Falls, Canada, 2008.
- [24] Zellig Harris. Distributional structure. *Word*, 10(2-3):1456-1162, 1954.
- [25] Thomas Hofmann. Unsupervised learning by Probabilistic Latent Semantic Analysis. *Machine Learning*, 42(1-2):177-196, 2001.

IJSER