



HAL
open science

Using Domain Ontologies for Classification and Semantic Interpretation of Documents

Samia Iltache, Catherine Comparot, Malik Si-Mohammed, Pierre-Jean Charrel

► **To cite this version:**

Samia Iltache, Catherine Comparot, Malik Si-Mohammed, Pierre-Jean Charrel. Using Domain Ontologies for Classification and Semantic Interpretation of Documents. International Workshop on Knowledge Extraction and Semantic Annotation (KESA 2016) in ALLDATA 2016: 2nd International Conference on Big Data, Small Data, Linked Data and Open Data, Feb 2016, Lisbon, Portugal. pp. 76-81. hal-01535945

HAL Id: hal-01535945

<https://hal.science/hal-01535945>

Submitted on 9 Jun 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Open Archive TOULOUSE Archive Ouverte (OATAO)

OATAO is an open access repository that collects the work of Toulouse researchers and makes it freely available over the web where possible.

This is an author-deposited version published in : <http://oatao.univ-toulouse.fr/>
Eprints ID : 16920

The contribution was presented at KESA 2016 :
<http://www.iaia.org/conferences2016/KESA.html>

To cite this version : Iltache, Samia and Comparot, Catherine and Si-Mohammed, Malik and Charrel, Pierre-Jean *Using Domain Ontologies for Classification and Semantic Interpretation of Documents*. (2016) In: International Workshop on Knowledge Extraction and Semantic Annotation (KESA 2016) in ALLDATA 2016 : 2nd International Conference on Big Data, Small Data, Linked Data and Open Data, 21 February 2016 - 25 February 2016 (Lisbon, Portugal).

Any correspondence concerning this service should be sent to the repository administrator: staff-oatao@listes-diff.inp-toulouse.fr

Using Domain Ontologies for Classification and Semantic Interpretation of Documents

Samia Iltache

UMMTO

Tizi Ouzou, Algeria

e-mail: s_iltache@hotmail.com

Catherine Comparot

IRIT, Université de Toulouse,
CNRS, INPT, UPS, UT1, UT2J
France

e-mail: Catherine.Comparot@irit.fr

Malik Si Mohammed

UMMTO

Tizi Ouzou, Algeria

e-mail: m_si_mohammed@esi.dz

Pierre-Jean Charrel

IRIT, Université de Toulouse,
CNRS, INPT, UPS, UT1, UT2J
France

e-mail: Charrel@univ-tlse2.fr

Abstract—The work presented in this paper addresses the problem of interpretation and semantic classification of documents. One of the issues faced by natural languages is related to the presence, in glossaries, of words with similar morphologies and different meanings. Our approach is based on the use of domain ontologies for nouns disambiguation. We begin our process with a global disambiguation, by linking the considered document to a semantic domain (represented by an ontology) which we select among several candidate ones. We define a candidate domain as any domain in which at least one significant word of the text can be considered and makes sense. We then perform a local disambiguation by using the selected ontology and finally build a semantic representation of the content of the document as a conceptual graph.

Keywords-Domain ontology; semantic interpretation; disambiguation; classification; conceptual graph.

I. INTRODUCTION

A document is represented by a set of words that expresses its global meaning. In conventional approaches, a document is represented by the lemmas of words describing its contents. To these lemmas is assigned a weight indicating their importance in the document. This weight combines local weighting linked to the document itself and a global weighting based on the considered corpus.

Semantic approaches aim to give meaning to the terms of the document to address the shortcomings of conventional indexing based on single words.

The issue of words with similar morphologies and different meanings is faced in all languages. If the assignment of adequate meaning to a word is easily done by a human being, because he uses his knowledge, this process is made difficult for an application using the textual content of the documents based on the morphological appearance of words.

Our approach aims to achieve an interpretation and semantic classification of textual content of documents.

We propose to use the knowledge represented by domain ontologies as a basis for our process. In fact, we consider that concepts of an ontology can allow to give the appropriate meaning to the words of the document. We first perform a global disambiguation through a classification process. This is to determine, among several domain ontologies, which one is the best to be considered, in order to obtain the correct semantic of document content; it is determined by the overall context of the document. In the second step, a local disambiguation is performed if some of the terms can be associated to several concepts within the retained ontology. The process, thus defined, allows to respond to the problem of polysemy and synonymy.

Our approach allows to thematically group the documents and to obtain a semantic representation of their content. A document can then be represented by a conceptual graph extracted from the ontology to which the document has been attached.

This paper is organized as follows. First, in Section II and Section III, we present a brief state of the art by introducing some works related to our problem. Then, in Section IV, we focus on the different steps of our process. In Section V, we present some examples to illustrate our approach. In the last Section, we conclude on the usefulness of our approach and give the prospects for its use and evolution.

II. SEMANTIC INDEXING, CONCEPTUAL INDEXING

To represent the meaning conveyed by the textual content of a document, several approaches use thesauri or ontologies to annotate the document. The semantic annotation is usually accompanied by a disambiguation process.

In order to find the appropriate meaning of an ambiguous word occurrence, endogenous approaches use its context in the document and all the documents of the corpus [1]. Exogenous approaches exploit external linguistic resources such as digital dictionaries or Machine

Readable Dictionary (MRD) [2], thesauri [3], or ontologies [4].

WordNet [5] is a linguistic resource. Its lexical database covers almost the entire English language. A concept in WordNet, which is called a synset, is represented by a set of synonyms. Synsets are connected by hyponym - hypernym (specialization - generalization) relations and meronymy - holonymy (part - all) relations. WordNet is a widely used resource, particularly in information retrieval. To represent a document, Baziz [6] defines a semantic core. The semantic content of a document is obtained by projecting the terms of the document on WordNet to extract the most representative synsets. The links between these synsets are weighted based on the semantic proximity (semantic similarity) between these synsets. The choice of synsets is based on two criteria: the co-occurrence called *cf.idf* and the semantic similarity used to disambiguate the synsets. Kolte [7] also uses WordNet to find the synsets corresponding to content of a document. He uses the various relationships defined in WordNet, as well as links, such as "ability link", "function link" and "capability link" to disambiguate the ambiguous words. For each word or group of words in a document d , Wang [8] constructs a matrix Uc for each candidate synset c , extracted from WordNet, corresponding to d . The rows and columns of Uc represent the words, di ($i=1,n$) forming c . The row i of Uc gives the probability that a word di and a word dj , ($j=1,n$) appear simultaneously in d . The matrix Uc denotes the relevance of d with a synset c . In WordNet, domains are assigned to synsets to define the different meanings they may have. Kolte [9] uses these domains to find the correct meaning of a synset depending on other terms appearing together with it in the same sentence. Fauceglia [10] disambiguates verbs by exploiting information about the verbs that appear in similar contexts. His approach is applied in the Event Mention Detection task (EMD) to classify event types. He uses a database of the meaning of the verbs and no structure highlighting a relationship between the meanings of the verbs is used as it is the case in WordNet.

III. AUTOMATIC CLASSIFICATION OF DOCUMENTS

The automatic text classification aims to organize documents into categories. One or more labels (classes, categories) are thus assigned to a document according to its text content.

Approaches dealing with supervised classification assign documents to predefined classes [11][12][13] while unsupervised classification approaches automatically define classes, called clusters [14].

In supervised classification, classifiers use two collections of documents: A collection containing learning documents to determine the features (terms) for each category and a collection containing new documents to be automatically classified. The classification of a new document depends on features retained for each category. A document is represented by a vector whose dimension is equal to the number of features selected to represent the different categories and no relationship between these features is highlighted. The vector document is then represented as a "bag of words".

Some classifiers create a "prototype" class from the learning collection [11]. This class is represented by the

average vector of all vectors of the documents in the collection. Only certain features are retained, which represents a loss of information.

Other approaches replace the learning collection composed of selected documents for each category, by data extracted from the "world knowledge" as Open Directory Project (ODP) [15]. Other approaches use thesauri [16] and domain ontologies [17] with conventional classifiers such as Support Vector Machine (SVM), Naïve Bayes, K-means, etc.) and represent a document by a vector of features represented by concepts or by a combination of terms and concepts.

The representation of the features by a vector assumes their independence from one another. The different approaches face the problem caused by the large size of the document vector, which reduces their performance. A step for restricting the features is thus performed.

IV. PROPOSED APPROACH

Our approach aims to build, for a document, a graph whose nodes and arcs are respectively represented by concepts and relations between concepts. Our process is based on a global disambiguation step based on a classification of documents using several domain ontologies and a local disambiguation based on a domain ontology.

The documents classification allows grouping documents according to the knowledge domain defined by their content. This grouping identifies a global similarity expressed by the context in which the document has a coherent sense. This classification determines the concepts to retain for the document through its global context.

The classification that we implement is a semantic classification because unlike conventional approaches, we take into account the link between terms with their context of appearance in the document and we extract concepts corresponding to these terms from domain ontologies.

The classification allows to project the content of a document on several domain ontologies to determine which one best expresses its content. Synonyms and polysemic terms are assigned to concepts representing their appropriate sense. We consider the following facts:

- Someone can use the same terms to describe different knowledge. Thus, a term may have several meanings depending on the context in which it is used. The same term ti extracted from a document d can then be assigned to several concepts which belong to different ontologies.

$$t_i^d = \{c_{\theta_1}, c_{\theta_2}, \dots\} \quad (1)$$

C_{θ_i} represents the concept extracted from the ontology θ_i .

- A term can match with several concepts of the same ontology.
- The theme discussed in a document depends on the terms used in its content and the way these terms are grouped together in sentences and paragraphs.

A. Projection, extraction of terms and candidate concepts.

The “projection” of a document on different ontologies allows to associate meaning to the terms of the document with respect to concepts belonging to these ontologies, and to select the candidate concepts. The notion of concept gives a meaning to a term relative to the domain in which this concept is defined.

We divide the whole document into sentences. Each sentence is browsed from left to right from the first word. We project the words of each sentence on different domain ontologies to extract the longer phrases (groups of words called "terms") that denote concepts. This choice is determined by: 1) the concepts are often represented by labels consisting of several words, 2) long terms are less ambiguous.

Several concepts belonging to the same domain ontology may be candidates for a given term.

B. Local disambiguation.

The disambiguation process is used to select for a term t the most appropriate concept among several candidates belonging to the same ontology. To do this, we consider the context of occurrence of the term t in the document.

We consider the following assumptions:

- We assume that the semantic link between the terms depends on the distance between these terms within the document. The shorter the distance, the greater the semantic link. The semantic link decreases when passing from sentence to paragraph and also from one paragraph to another.
- We choose the appropriate concept for the term t , taking into account both the semantic distance between the term t with neighboring terms, (i.e. which occur in its context), and the semantic distance between concepts associated with the term t and the concepts corresponding to the neighboring terms in the ontology considered.
- The meaning of a term t in a document is determined by its nearest neighbors terms. t will then be disambiguated by its nearest neighbor on the left or by its nearest neighbor on the right. In case the left and right neighbors exist simultaneously, they will both be taken into consideration.

The disambiguation process is then done in three levels, starting at the sentence level. For each sentence, the ambiguous terms are disambiguated considering their left and right neighbors in the sentence. Any disambiguated term helps to move forward in the process of disambiguation of next terms. This process is repeated in case ambiguous terms still remain, considering in a second step the paragraph level, and finally, if necessary, the document level.

The disambiguation of a term t at sentence level is represented in Figure 1.

The disambiguation process at sentence level considers neighboring terms, unambiguous, that have associated concepts in the ontology considered, surrounding t : it retrieves Cv_g and Cv_d , corresponding respectively to v_g , the nearest neighbor on the left of t and v_d , the nearest neighbor on the right of t .

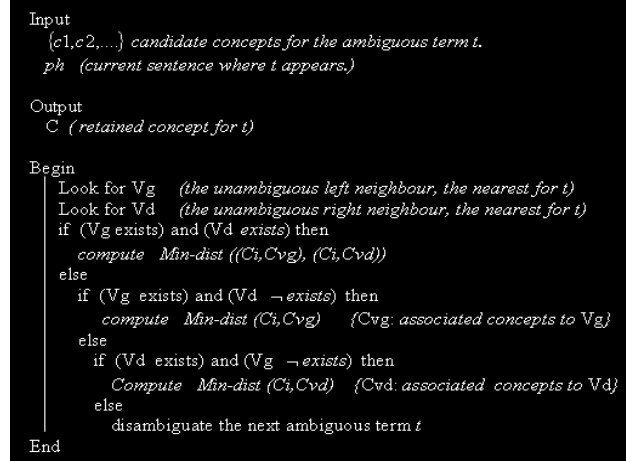


Figure. 1. Local disambiguation process, sentence level.

The appropriate concept for the term t among candidate concepts is the semantically nearest concept of Cv_g or Cv_d . This amounts to browsing the ontology and calculating the minimum distance between each concept associated with t and candidate concepts Cv_g , Cv_d . Several existing metrics in the literature are used to calculate this minimum distance.

C. Classification: global disambiguation

While Kolt [9] determines the meaning of an ambiguous word with the most represented domain identified by the terms appearing with it in the same sentence, we seek to determine the context defined by a document. We propose to represent it not by words but by a set of concepts.

We rely on Wang's approach [8] which operates the occurrence of words within paragraphs to determine which concept to assign to a term of a document. We extend the process to the classification in order to determine the importance of all concepts extracted from different ontologies relative to the terms of the document.

At the end of the preceding steps, a document d is represented by several set of concepts extracted from domain ontologies θ_i on which it has been projected.

$$d = \begin{cases} \theta_1^d = \{c_{11}, c_{21}, \dots, c_{n1}\} \\ \theta_i^d = \{c_{1i}, c_{2i}, \dots, c_{ni}\} \\ \dots \end{cases} \quad (2)$$

The classifier needs to conclude the relevance of a document relative to a given context and to choose among the different ontological representations, which one best corresponds to its context. To do this, associating different domain ontologies to classes, the classifier will make the classification of a document relative to a single domain ontology.

The words used to describe a particular idea are not arbitrarily chosen. They are semantically related and are chosen with a common sense guided by this idea. However, it is almost impossible to find a document or a

text in which all used terms refer exclusively to a same domain.

Recall that the previous steps are used to extract the concepts corresponding to the terms in the document. The extracted concepts can be related to multiple ontologies.

The classification we define in this work aims to determine, for each ontology θ_i , the semantic weight of each concept extracted for the document d . This determines the importance of a concept relative to a document. The evaluation of this weight is performed at two levels: paragraph level and document level.

Paragraph level: We calculate the weight of each concept C_i based on the other concepts appearing with it in a paragraph.

Document level: We calculate the total weight of each concept C_i throughout the document. This weight is obtained by adding the weights obtained for the concept C_i in the various paragraphs of the document d .

For each ontology and for each document we associate a matrix such as (3).

$$M_{\theta_i}^d = \begin{pmatrix} lc_1c_1 & lc_1c_2 \dots & lc_1c_n \\ \vdots & \vdots & \vdots \\ lc_nc_1 & lc_nc_2 \dots & lc_nc_n \end{pmatrix} \quad (3)$$

The rows and columns of this matrix represent all concepts extracted from ontology θ_i for the document d .

C_i is any concept extracted from the ontology θ_i after the projection of the document d on θ_i ; $lcicj$ represents the weight of the link between the concept C_i and the concept C_j ($i \neq j$). This weight is calculated as follows:

- The matrix is initialized to zero
- If a term ti and a term tj appear together within the same paragraph of the document d and concepts C_i and C_j correspond to terms ti and tj respectively, then the weight $lcicj = 1$.
- The weight $lcicj$ is updated each time terms ti and tj appear together in the same paragraph.
- The weight $lcicj$ corresponds to the appearance of term ti in the paragraph. It is equal to 1.
- The weight $lcicj$ is updated for all paragraphs in the document d .

Each row of the matrix represents the total weight of a concept extracted from the ontology θ_i relative to a document d . This weight assesses the importance of the concept C_i in d .

The total weight of all the extracted concepts of an ontology relative to document d , measures how well each ontology represents this document. The highest score will determine the ontology candidate which will be chosen to represent the document d .

V. IMPLEMENTATION AND EXAMPLES

We implemented our approach using both WordNet and WordNet Domains resources. In WordNet Domains, several knowledge domains are used, such as medicine, computer science, economy etc, and each synset is

annotated with one or more domains in which it has a meaning.

To achieve our classification, we have assimilated these domains to domain ontologies. To evaluate the distance between two synsets in WordNet we used Rita similarity metric [18].

The words within sentences are tagged with their type (noun, verb, adverb, adjective, etc.) by Stanford Part-Of-Speech Tagger (POS Tagger) [19].

To illustrate our approach, we apply it on the three following examples:

- Txt1: *The role of banks in the economy was clear and well established as the financial markets were underdeveloped because they were the only ones to provide liquidity and credit to businesses and households. The unprecedented development of financial markets, driven by the late 1970s in the Anglo-Saxon countries, has led some economists to question the specificity of bank financing compared with direct funding and the survival of traditional banks. Several arguments have been advanced.*
- S1: *Banks use their networks to exploit economies of scale between activities (collection of savings, management of means of payment, exchange, offer insurance products, securities investment services....*
- S2: *The player throws the baseball and he improves the score...*

A. Example 1: the Txt1 case

1) *Global disambiguation:* We consider four domains and we apply the classification process to determine the domain that represents best the content of the text *Txt1*. It determines the synsets to retain for the text through its global context. Table I shows the result of the projection of *Txt1* on four ontologies and the score obtained by each ontology.

The selected domain is *Economy* because it has obtained the highest score.

2) *Local disambiguation, sentence level:* In the domain retained, the term *economy* has two synsets. So this is an ambiguous term. A local disambiguation is performed to determine what synset to retain for this term. It is performed at sentence level. The nearest unambiguous neighbor of the term *economy* is only on the right: it is the term *credit*. There is no path between *economy* and *credit* in WordNet. Another nearest neighbor is sought in the sentence. This is the term *business* that is on the right of *economy*.

The distance between *business* 07485368-n and *economy* 00182005-n is: 1.0.

The distance between *business* 07485368-n and *economy* 07857433-n is: 0.8333333.

The synset retained for *economy* is 07857433-n.

The text *Txt1* is represented by the following synsets (economy 07857433-n, credit 12616435-n, business 07485368-n, economist 09401295-n).

TABLE I. BREAKDOWN BY ONTOLOGIES OF SYNSETS ASSOCIATED TO TERMS OF TXT1.

Ontologies	Terms	Synsets	Score	
economy	economy	00182005-n 07857433-n	16	
	credit	12616435-n		
	business	07485368-n		
	economist	09401295-n		
finance	bank	12599211-n	1	
enterprise	business	01033295-n 01031794-n 07571175-n	6	
		financing		01036077-n
		funding		01036077-n
banking	bank	02690337-n 07909067-n	8	
		credit		12620638-n
Synsets	Definitions (Glosses of WordNet)			
00182005-n	an act of economizing; reduction in cost.			
07857433-n	the system of production and distribution and consumption.			
07485368-n	business concerns collectively. "Government and business could not agree"			
01033295-n	the volume of business activity; "business is good today"			
01031794-n	commercial_enterprise, business_enterprise the activity of providing goods and services involving financial and commercial and industrial aspects.			
07571175-n	business_organisation a commercial or industrial enterprise and the people who constitute it.			

B. Example 2: the S1 case

S1 is an extract of a sentence belonging to a text classified in the domain *Economy*. Table II summarizes the synsets associated with its terms.

Payment is an ambiguous term since it has two synsets. A local disambiguation is realized at sentence level. The nearest unambiguous neighbors for *payment* are *means*, which is on the left and *exchange* which is on the right.

The distance between *exchange* 01045967-n and *payment* 01056649-n is: 0.4.

The distance between *exchange* 01045967-n and *payment* 12522505-n is: 1.0.

The distance between *means* 12596703-n and *payment* 01056649-n is: 1.0.

The distance between *means* 12596703-n and *payment* 12522505-n is: 0.85714287.

The shortest distance is given by the term *exchange* that is on the right of *payment*. The synset retained for *payment* is 01056649-n.

C. Example 3: the S2 case

We consider an extract from the sentence S2 belonging to a text classified in the domain *Play*. Table III summarizes the synsets associated with its terms.

Baseball is ambiguous. It has two neighbors but only one is unambiguous. This is the term *player* that is on the left. So *Player* disambiguates *baseball*.

The distance between *player* 09762180-n and *baseball* 02701461-n is: 0.75.

TABLE II. SYNSETS ASSOCIATED TO TERMS OF S1

Terms	Synsets	Definitions (Glosses of WordNet)
economy of scale	00182453-n	
saving	00182005-n	
means	12596703-n	
payment	01056649-n	the act of paying money.
	12522505-n	a sum of money paid.
exchange	01045967-n	
security	12592487-n	
investment	12576508-n	

TABLE III. SYNSETS ASSOCIATED TO TERMS OF S2

Terms	Synsets	Definitions (Glosses of WordNet)
player	09762180-n	
baseball	02701461-n	a ball used in playing baseball.
	00447188-n	a ball game played with a bat and ball between two teams of 9 players.
score	00176295-n	
	12829162-n	

The distance between *player* 09762180-n and *baseball* 00447188-n is: 1.0

The synset retained for *baseball* is: 02701461-n.

VI. CONCLUSION

In this paper, we proposed an approach to extract semantics from documents by using domain ontologies with a disambiguation process. This process combines local and global disambiguation. The first one finds an appropriate concept for a term with several meanings in a single domain ontology, the second one retrieves the appropriate concept for a term that has multiple meanings in different knowledge domains. Throughout the disambiguation process, we took into account the context of appearance of the ambiguous terms in the document. The quality of the disambiguation process of course depends on domain ontologies, since they must cover the entire vocabulary of the represented domain.

The major problem conventional classifiers suffer from is the vector representation of a document in a high dimensional space. Indeed, the size of the vector equals the number of features that represent all classes used by the classifier. This dimension, very large, lowers the performance of classifiers. Moreover, characteristics representing a document are independent of each other.

Our approach has the advantage of responding to the problem of polysemy and synonymy engendered by the terms of the document. The document is not represented by a vector of high dimension but by a conceptual graph where concepts correspond only to the terms describing its contents. We believe that the use of ontologies in our classification process is a more stable base than the use of a set of learning documents, in which the choice of such learning documents affects the result of the classification.

We have conducted tests on a first set of short documents. Even if the obtained results are very encouraging, we have obviously to confirm the interest of our approach by considering larger collections, with more candidate domains. This is what we plan to do in order to

assess the effectiveness of our approach in comparison to the existing ones.

REFERENCES

- [1] H. Schütze and J. Pedersen, "Information retrieval based on word senses," In Proceedings of the 4th Annual Symposium on Document Analysis and Information Retrieval, 1995, pp. 161-175.
- [2] J.A. Guthrie, L. Guthrie, Y. Wilks and H. Aidinejad, "Subject-dependant cooccurrence and word sense disambiguation," In Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics, Berkley, CA, 1991, pp.146-152.
- [3] D. Yarowsky, "Word-Sense Disambiguation Using Statistical Models of Roget's Categories Trained on Large Corpora," Proceedings of the 14th international Conference on Computational Linguistics (COLING-92). Nantes, France, August. 1992, pp.454-460.
- [4] P. Resnik, "Disambiguating noun groupings with respect to WordNet senses," 3th Workshop on Very Large Corpora, 1995, pp.54-68.
- [5] G. A. Miller, "WordNet: A Lexical Database for English," Communications of the ACM Vol. 38, No. 11, 1995, pp. 39-41.
- [6] M. Baziz, M. Boughanem and N. Aussenac-Gilles, "Conceptual Indexing Based on Document Content Representation," CoLIS 2005, pp. 171-186.
- [7] S. G. Kolte and S. G. Bhirud, "Exploiting links in WordNet hierarchy for word sense disambiguation of nouns," International Conference on Advances in Computing, Communication and Control, (ICAC3'09), 2009.
- [8] H. Wang, Y. Guo and X. Shi, "Research of the conceptual representing of documents based on light ontology," 9th International Conference on Fuzzy Systems and Knowledge Discovery, (FSKD, 2012), 2012.
- [9] S. G. Kolte and S. G. Bhirud, "WordNet: A Knowledge Source for Word Sense Disambiguation," International Journal of Recent Trends in Engineering, Vol 2, No. 4, November. 2009.
- [10] N. R. Fauceglia, Y.C. Lin, X. Ma and E. Hovy, "Word sense disambiguation via propstore and ontonotes for event mention detection," In Proceedings of the The 3rd Workshop on EVENTS: Definition, Detection, Coreference, and Representation, Denver, Colorado, June. 2015, pp.11-15.
- [11] T. Joachims, "A Probabilistic Analysis of the Rocchio Algorithm with TFIDF for Text Categorization," Proceedings of ICML-97, Tennessee, 1997, pp.143-151
- [12] Y. Yang and X. Liu, "A re-examination of text categorization methods," 22nd Annual International SIGIR, Berkley, August. 1999, p. 42-49.
- [13] S. Jaillet, A. Laurent and M. Teisseire, "Sequential patterns for text categorization". Intelligent Data Analysis, IOS Press, 2006.
- [14] A. Hotho, A. Maedche and S. Staab, "Ontology-based Text Document Clustering," KI 16 (4), 2002, pp. 48-54.
- [15] E. Gabrilovich and S. Markovitch, "Feature Generation for Text categorization Using World Knowledge," IJCAI 2005: the Nineteenth International Joint Conference on Artificial Intelligence, Edinburgh, Scotland, UK, July 30-August 5.2005, pp. 1048-1053
- [16] A. Hotho, S. Staab and G. Stumme, "Ontologies Improve Text Document Clustering," ICDM:3rd IEEE International Conference on Data Minin 2003, pp. 541-544
- [17] H. H. Tar and T.T. Soe.Nyunt, "Ontology-Based Concept Weighting for Text documents," International Conference on Information Communication and Management IPCSIT vol.16, IACSIT Press, Singapore, 2011.
- [18] D. C. Howe, "RiT: creativity support for computational literature," In Proceedings of the seventh ACM conference on Creativity and cognition (C&C '09). ACM, New York, NY, USA, 2009, pp. 205-210.
- [19] K. Toutanova, D. Klein, C. Manning, and Y. Singer, "Feature-Rich Part-of-Speech Tagging with a Cyclic Dependency Network", In Proceedings of HLT-NAACL, 2003, pp. 252-259.