

Intelligent text classification system based on self-administered ontology

Manoj MANUJA^{1,2,*}, Deepak GARG¹

¹Department of Computer Science and Engineering, Thapar University, Patiala, India

²Education and Research Department, Infosys Ltd., Chandigarh, India

Received: 14.05.2013

Accepted/Published Online: 03.09.2013

Printed: 28.08.2015

Abstract: Over the last couple of decades, web classification has gradually transitioned from a syntax- to semantic-centered approach that classifies the text based on domain ontologies. These ontologies are either built manually or populated automatically using machine learning techniques. A prerequisite condition to build such systems is the availability of ontology, which may be either full-fledged domain ontology or a seed ontology that can be enriched automatically. This is a dependency condition for any given semantics-based text classification system. We share the details of a proof of concept of a web classification system that is self-governed in terms of ontology population and does not require any prebuilt ontology, neither full-fledged nor seed. It starts from a user query, builds a seed ontology from it, and automatically enriches it by extracting concepts from the downloaded documents only. The evaluated parameters like precision (85%), accuracy (86%), AUC (convex), and MCC (high positive) demonstrate the better performance of the proposed system when compared with similar automated text classification systems.

Key words: Ontology, support vector machine, resource description framework, text classification

1. Introduction

The World Wide Web is the biggest and most current source of information on any and every domain on this earth. This set of information may contain documents ranging from best practices to technical reports, customer feedback, and product review comments, to name just a few. Around 80% of this information is written in natural language and unstructured in nature [1]. Often a novice user finds it very difficult to search for useful information on the web without prior knowledge of the subject supported by rich clues. This is primarily because of the fact that web classification systems (mainly search engines) respond to any user query on the basis of its syntax instead of searching on the basis of its semantics.

Computing for human experience (CHE) is one of the futuristic thoughts that provide a vision for future man-machine interfaces with a realistic implementation [2]. CHE talks about technology-rich intelligent systems enabling human experience to gather and apply knowledge in relevant fields of sentiment and opinion-mining. The CHE vision motivates authors to design and develop an intelligent system that can help to classify data available on the web with minimal explicit effort being put forth by humans. In today's world of the semantic web, it is indeed relevant to have any intelligent classification system be based on semantics rather than syntaxes.

Much work has been done on developing semantic-based classification systems during the past few years [3-5]. These systems make all forms of information be linked through semantics so that humans can utilize this wealth of information automatically using machine learning methods and algorithms. Rapid growth of

*Correspondence: manojmanuja2670@gmail.com

linked data in the recent past has led to the availability of domain-specific ontologies on the web in abundance [6,7]. Ontologies provide a controlled vocabulary of concepts along with their relations explicitly defined using machine processable semantics. It is very time-consuming for humans to build a machine-understandable relationship graph covering all the domains available around us. Hence, it is imperative to automate the ontology learning process, which in turn helps in enhancing human-machine interactions [8]. Quite a few frameworks are available that address this automation problem, as mentioned in the next section where related work in this field is discussed. Most of them use full-fledged prebuilt domains or seed ontology to start the automation and enrichment process.

In this paper, we share the details of a proof of concept (PoC) of a classification system that is self-governed in terms of ontology population and does not require any prebuilt ontology, neither full-fledged nor seed. It starts from a user query, builds a seed ontology from it, and automatically enriches it by extracting concepts from the downloaded web documents only. The suggested system framework facilitates the human-machine interaction in line with the relevant search results based upon the user query and classifies the downloaded documents in appropriate categories. The important point of our framework is that it does not use any prebuilt ontology and starts from scratch to convert the knowledge into a powerful web document classification mechanism purely focused on the user's query.

The paper is organized as follow. Section 2 provides the related work in this field. Section 3 provides a complete view of the system framework. Section 4 gives insight onto the experimental results achieved during the implementation, with Section 5 providing a detailed analysis. Comparison of our framework with the similar closest frameworks is done in Section 6. Section 7 concludes the paper.

2. Related work

There are many frameworks that have been suggested and implemented by various researchers that describe automatic ontology learning for the semantic web.

Maedche and Staab [9] suggested a comprehensive 5-step framework of ontology learning for the semantic web. The framework proceeds through importing, extracting, pruning, refining, and evaluating the ontology. They used the Text-To-Onto ontology learning environment, which helps in learning from free text, dictionaries, or legacy ontologies to build and enrich a given domain ontology. The suggested framework is semiautomatic as it requires the involvement of an ontology engineer to support the framework during different stages of learning.

Navigli et al. [10] suggested the OntoLearn system for automatic ontology learning from domain text. This system crawls domain-specific web sites and data warehouses for extracting terminologies, and then filters them using natural language processing and statistical techniques. A domain concept forest is created that provides a semantic interpretation of these terminologies duly supported by WordNet and SemCor lexical knowledge bases. The framework has three major phases, namely terminology extraction, semantic interpretation, and creation of the WordNet specialized view. The inductive machine learning technique is used to associate the appropriate relations among complex components of the domain concept. Again, this framework is semiautomatic as it relies on WordNet and requires the involvement of ontology engineer.

Wei et al. [11] suggested an ontology-based approach to classify web documents. An already available knowledge base is used as a starting point. Ontology is built for each subclass of the knowledge base, which uses the Resource Description Framework Schema for transforming knowledge into ontology. A comparison between various machine learning algorithms like the support vector machine (SVM), k-nearest neighbor, and latent semantic association with the ontology-based approach clearly showed the advantages of the ontology-based classifier. Again, this system also depends on prior available information in the form of a knowledge base.

Luong et al. [12] suggested a framework for ontology learning that uses a web crawler to retrieve documents from the web, identifying domain-specific documents using a support vector machine and extracting useful information from them to enrich domain-specific ontology. Existing small, manually created domain ontology is being enriched through this process by adding new concepts in its hierarchical structure. Our approach is similar to this framework with a difference of eliminating the need of any external preexisting domain ontology to classify the new unknown documents.

Brank et al. [13] suggested a framework to classify multiclass textual documents using a SVM and a coding matrix. The existing ontology is populated by extracting a hierarchy of concepts and instances from the large corpus of documents. Again, the main assumption here is that some “training data” are already available that consist of primarily a set of instances with correct assignment of concepts.

Speretta and Gauch [14] also provided a semiautomatic approach to enrich the vocabulary of each concept in a given ontology with words mined from the set of crawled documents and then combined with WordNet.

3. System framework

The proposed system framework is triggered by the user input in the form of English query sentences. This moves to the query manager, which in turn provides the same to two different blocks, namely the focused web crawler and seed ontology generator. After downloading the document corpus from the web, preprocessing is carried out, which in turn becomes the input to feature extraction and selection blocks. The output of this block is again shared with two blocks in parallel, one being the ontology manager for ontology population and the other being the SVM classifier for training purpose. After the ontology is populated, it is also provided to the training model, which is replicated as a testing model to finally evaluate the model performance in terms of classifying documents based on the user query. The block-wise details of the system are provided in Figure 1.

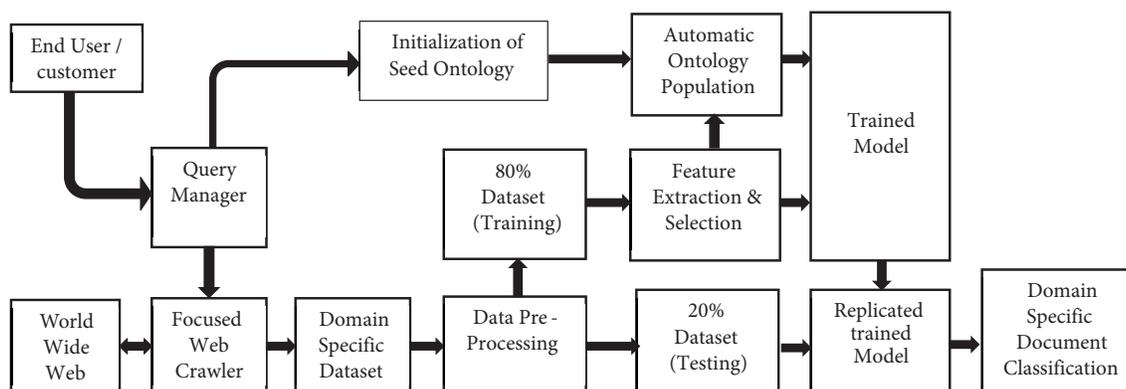


Figure 1. System design of self-governed ontology-based classification system.

End-user query interface: The query from the end-user is presented to the query manager, which is a GUI-based Interface. The query interface asks for certain information from the user. It tries to collect as much information as it can in terms of concepts being shared by the user in the form of unstructured English sentences.

Query manager: This interface handles two tasks. The first task is to make the query available to the focused web crawler, which searches relevant web pages from the web. The second task is to provide the basic information fed by the user to the block, which initiates the building of “Seed Ontology” from it.

Focused web crawler: It is used to retrieve documents and information from the web purely based on the query submitted by the user. The outcome of the focused web crawler is a corpus of web documents that

are labeled as the domain-specific dataset and stored locally. This is a set of raw web documents. Although focused towards the domain of the query shared by the user, the dataset still requires preprocessing to be done before we can focus on the information extraction process.

Data preprocessing: Preprocessing is done to eliminate language-dependent factors. This step is very critical and mandatory prior to doing any meaningful text mining or analytics. The basic steps are:

- Scope: Choose the scope of the text to be processed. In our case, it is a set of documents.
- Tokenization: Break text into discrete words called tokens.
- Remove stop-words: Remove common words such as ‘the’, ‘they’, etc.
- Normalize spellings: Unify misspellings and other spelling variations into a single token.
- Detect sentence boundaries: Mark the end of sentences.
- Normalize case: Convert the text to either all lowercase or all uppercase.
- Stemming: Remove prefixes and suffixes to normalize words. For example run, running, and runs would all be stemmed to ‘run’.

Here we can define feature extraction as a combination of tokenization, stop-word removal, and stemming. We have used term frequency-inverse document frequency (TF-IDF) to extract features in the corpus.

TF-IDF is an important technique in information retrieval that evaluates how important a word is in a document [15]. It also plays an important role in converting the textual representation of information into a vector space model or into sparse features. TF-IDF determines the relative frequency of the words in a specific document as compared to the inverse proportion of that word over the entire corpus of the documents under review.

Let D = Collection of documents, w = total number of terms in a document, t = a term, and d = and individual document where $d \in D$.

Then we have term-frequency defined as $tf(t, d) = \frac{f(t, d)}{\max\{f(w, d) : w \in d\}}$.

The IDF is defined as a measure of whether the term t is common or rare across all documents:

$$idf(t, D) = \log \frac{|D|}{|\{d \in D : t \in d\}|},$$

where $|D|$ is the total number of documents in the corpus and $|\{d \in D : t \in d\}|$ is number of documents where the term t appears, i.e. $tf(t, d) \neq 0$.

In the case that the term is not in the corpus, then the denominator will become “divide-by-zero”. Therefore, we adjust the formula as $1 + |\{d \in D : t \in d\}|$ to deal with this scenario.

Hence, $tf - idf$ is calculated as $tf - idf(t, d, D) = tf(t, d) \times idf(t, D)$.

In other words, $tf - idf$ assigns each term present in the document a weight that is:

- Highest when a term t occurs many times within a small number of documents,
- Lower when the term occurs fewer times in a given document or occurs in many documents,
- Lowest when the term occurs in virtually all documents.

After the dataset is preprocessed and ready for further processing, we split the whole set into two parts, whereby 80% of the dataset is used for training purposes and 20% is used for testing the trained model. This split is done randomly and no specific technique is used.

Feature selection: A minimal subset of features (extracted in the previous step) is selected so that we may realize the maximum generalization ability of the classifier. Two well-established methods are available in machine learning for feature selection, namely wrappers and filters [16]. Wrapper methods are very time-consuming and hence have been ignored in this exercise. Filter methods work independently of the learning algorithm that will use the selected features. During feature selection, the filter method uses an evaluation metric that measures the ability of the feature to differentiate each class from the other. There are two types of filter methods, namely forward selection and backward selection. In backward selection, all the features are considered in the first instance and one feature is deleted at a time, which deteriorates the selection criteria the least. We go on deleting the features until the time selection criteria reach a particular acceptable value. In forward selection, an empty set of features is the starting point. We go on adding one feature at a time, which improves the selection criteria the most.

Selected features during this step are used for two purposes. One is to generate the training model with the SVM as the learning algorithm. The second is to employ these selected features as inputs to the ontology manager, which populates the ontology in the context of the user query.

Seed ontology generator: This interface takes the preprocessed user query from the query manager as input. This input is converted into a basic ontology tree using a resource description framework (RDF) graph [17]. Any given RDF graph contains a collection of triples; each consists of a combination of subject – predicate – object (S – P – O). Each triple is extracted from a given sentence, which reflects the relationship between its subject and object linked by a predicate. A sample RDF graph is given below.

```
<?xml version = "1.0"? >
< rdf:RDF xmlns:rdf = "http://www.w3.org/1999/02/22-rdf-syntax-ns#" >
< rdf:Description about = "http://www.whitehouse.gov/~BarackObama/" >
.
.
< /rdf:Description >
< rdf:Description rdf:ID = "Barack Obama" >
.
.
< /rdf:Description >
< /rdf:RDF >
```

We term it as seed ontology equipped with a very small number of classes. This is the starting point of enriching the ontology for the specific search query written by the user.

Ontology manager: This is the most critical block of our whole system. The success of the system depends on how this block populates the ontologies in the form of RDF triples from the given set of documents downloaded from the web using the focused web crawler. The process flow for the ontology manager is shown in Figure 2.

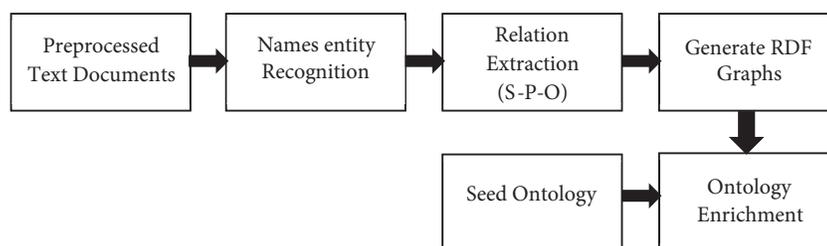


Figure 2. Process of ontology population by ontology manager.

The process starts with preprocessed text documents as input. Named entities are recognized and relations are extracted to primarily mark subjects, objects, and predicates followed by RDF translation of S – P – O. This block takes two inputs, one being the seed ontology prepared by the seed ontology generator and features extracted by the feature extraction block.

Training and testing model: The training model is built with two styles, one being built using machine learning algorithms and the second being built using the ontology prepared by the ontology manager. This model is then replicated as a testing model to check the accuracy and performance of the model using an unknown dataset. Document classification is done according to the user-based query and categorized into appropriate groups.

4. Experimental results

The proposed system is evaluated using two use-cases, the first being evaluated with three offline datasets and the other being evaluated using an online dynamic dataset downloaded from the web. We have done experiments with the below mentioned two setups using the LIBSVM package [18]:

- SVM-based classifier with linear kernel,
- SVM-based classifier with ontology based approach.

Ten-fold cross-validation (CV) is used during the experimentation stage. Ten-fold CV primarily breaks the given data into 10 sets of equally sized subsets, training on 9 datasets and testing on 1 dataset. The whole process is repeated 10 times and the accuracy of the classification process is calculated by taking the mean of all the stages. Ten-fold CV is a useful tool for accuracy estimation and model selection [19].

Performance metrics:

The following performance measures are evaluated:

P = number of relevant documents classified as relevant (true positive),

Q = number of relevant documents classified as not relevant (true negative),

R = number of not relevant documents classified as relevant (false negative),

S = number of not relevant documents classified as not relevant (false positive).

Therefore, the total number of documents $T = (P + Q + R + S)$.

The performance measure parameters are primarily precision, recall, and F-measure.

Precision = number of correctly identified items as percentage of number of items identified = $P/(P+S)$.

Recall = number of correctly identified items as percentage of total number of correct items = $P/(P+R)$.

F-measure = weighted average of precision and recall = $(2 \times \text{precision} \times \text{recall}) / (\text{precision} + \text{recall})$.

Accuracy = degree of conformity of a measured quantity to its true value = $(P+Q)/T$.

Accuracy is the degree of conformity of a measured quantity to its true value, while precision is the degree to which further measurements show similar results. In other words, the precision of an experiment is a measure of the reliability of the experiment, whereas the accuracy of an experiment is a measure of how closely the experimental results agree with a true value. In our case, we have compared both accuracy and precision parameters to provide a holistic insight into the experimental outcomes under different use-cases.

The Matthews correlation coefficient (MCC) is a correlation coefficient between the observed and predictive binary classification, returning a value between +1 and -1 with +1 providing perfect prediction and -1 providing total disagreement [20].

$$\text{MCC} = (P \times Q - S \times R) / \sqrt{((P + S)(P + Q)(S + R)(Q + R))}.$$

True positive rate (TPR) = number of true positives divided by the total number of positives = $P/(P+S)$.

False positive rate (FPR) = number of false positives divided by the total number of negatives = $S/(Q+R)$.

The area under the receiver operator curve (AUC) depicts a trade-off between benefits of the system (i.e. true positives) and cost overhead to the system (i.e. false positives). The receiver operator curve (ROC) is drawn with the TPR on the y-axis and the FPR on the x-axis.

Use-Case1: Offline classification:

Three benchmark datasets, namely Reuters-21578 [21] (dataset available at <http://www.daviddlewis.com/resources/testcollections/reuters21578/>), 20 Newsgroups collection [22] (dataset available at <http://qwone.com/~jason/20Newsgroups/>), and WebKB [23] (dataset available at <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-20/www/data/>), are used for offline classification. Reuters-21578 is currently the most widely used test collection for text categorization research. We have 6532 documents as the training dataset and 2568 documents as the test dataset. The 20 Newsgroups dataset is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups. The data are organized into 20 different newsgroups, each corresponding to a different topic. We have segregated the whole dataset into two parts with 11,293 documents as the training dataset and 7528 documents as the test dataset. WebKB is a collection of web documents collected by the World Wide Knowledge Base (Web -> Kb) project of the CMU text learning group, downloaded from the 4 Universities Data Set Homepage. These pages were collected from computer science departments of various universities in 1997, manually classified into four main classes: student, faculty, course, and project. The dataset is divided into two parts with 2803 documents forming the training dataset and 2396 documents forming the test dataset.

It is clear from Table 1 that the SVM with 10-fold CV appears to be better than the ontology-based classifier in some scenarios where the user query is straight and simple. As we go on adding the complexity to the user query, our proposed model built on ontology performs comparatively well. Average F-measure (85% for SVM with system-built ontology, 84% for SVM only) and accuracy (86% for SVM with system-built ontology, 77% for SVM only) parameters are much better in case of the SVM with system-built ontology framework as compared to the simple SVM-based system.

Use-Case2: Online classification:

We have done the experiment with 5 different user query strings. The focused web crawler downloads the top 40 pages each from five search engines, namely Google, Bing, Yahoo, AltaVista, and AOL, and 160 pages were preprocessed and then used for feature extraction and ontology population. The remaining 40 pages were thrown to the testing model for run time classification. Table 2 provides the top 5 categories of documents for each of the user query strings.

Table 1. Performance metrics and accuracy of 3 benchmark DBs under User-Case1.

Dataset	User query string	SVM only			SVM + system-built ontology				
		Precision	Recall	F-measure	Accuracy	Precision	Recall	F-measure	Accuracy
Reuters-21578	American-Samoa	0.91	0.89	0.9	0.88	0.9	0.88	0.89	0.92
	Association of International Bond Dealers	0.84	0.86	0.85	0.75	0.86	0.87	0.86	0.89
	African Development Bank	0.87	0.89	0.88	0.84	0.86	0.85	0.85	0.78
	Lexan Polish	0.88	0.9	0.89	0.81	0.86	0.85	0.85	0.88
20-Newsgroups	NASA	0.92	0.91	0.91	0.86	0.9	0.91	0.9	0.92
	What is a squid?	0.85	0.84	0.84	0.74	0.87	0.89	0.88	0.81
	Course at Cornell	0.81	0.79	0.8	0.72	0.85	0.84	0.84	0.74
WebKB	Raymond J. Mooney	0.76	0.78	0.77	0.65	0.8	0.79	0.79	0.91
	Internet Softbot	0.77	0.75	0.76	0.66	0.78	0.79	0.78	0.88

Table 2. User Query response under Use-Case2.

User query string	Top Category - 1	Top Category - 2	Top Category - 3	Top Category - 4	Top Category - 5
Barack Obama	Wikipedia	Personal Web Site	White House	Biography	News
US President Barack Obama	White House	Wikipedia	Personal Web Site	News	Biography
Barack Obama holiday in Hawaii	News	Wikipedia	-	-	-
Barack Obama visited China	News	Wikipedia	-	-	-
Sandy	Sandy Hurricane-News	Sandy Hook Schoo-News	Sandy - City	Sandy - Wikipedia	Sandy - social networking sites

It is very clear that “Wikipedia” is the top classified set of documents when the query string says “Barack Obama”. It changes to “White House” when the query string is modified to “US President Barack Obama”. Queries 3 and 4 also show that when the user is trying to search for some happening, the top category of classified documents changes to “News”. The fifth query gives some uneven results because of the query string meaning. It is clear from the top 5 categories that although we got related “News” under the top 2 categories, the classification system provides an option to the user to choose the classification folder accordingly to his/her choice. This particular output is encouraging for us where we have provided the novice user an opportunity to choose from these five categories when “Sandy” is searched.

The experimental results are quite heartening. It is clear from Table 3 that the SVM with system-built ontology classifier is able to provide slightly better results than the simple SVM classifier. Although average accuracy is improved by a factor of 0.3, average F-measure improvement is just 0.1.

Table 3. Performance metrics under Use-Case2.

User query string	SVM only				SVM + system-built ontology			
	Precision	Recall	F-measure	Accuracy	Precision	Recall	F-measure	Accuracy
Barack Obama	0.84	0.85	0.84	0.85	0.83	0.84	0.83	0.89
US President Barack Obama	0.85	0.88	0.86	0.82	0.85	0.83	0.84	0.91
Barack Obama holiday in Hawaii	0.84	0.85	0.84	0.82	0.85	0.86	0.85	0.82
Barack Obama visited China	0.85	0.83	0.84	0.83	0.86	0.84	0.85	0.83
Sandy	0.8	0.79	0.79	0.81	0.82	0.81	0.81	0.82

5. Analysis

Results of the ontology-based approach for both user-cases, i.e. with offline datasets and live datasets from the web, reflect that the proposed system performs reasonably well while classifying documents. The comparison of average performance and accuracy parameters for both classifiers (SVM and SVM with system-built ontology) is shown in Figure 3. Average MCC values calculated from various performance parameters are shown in Table 4. All positive values of MCC (>0) depict a reasonable representation of quality predictions as received from the experimental setup.

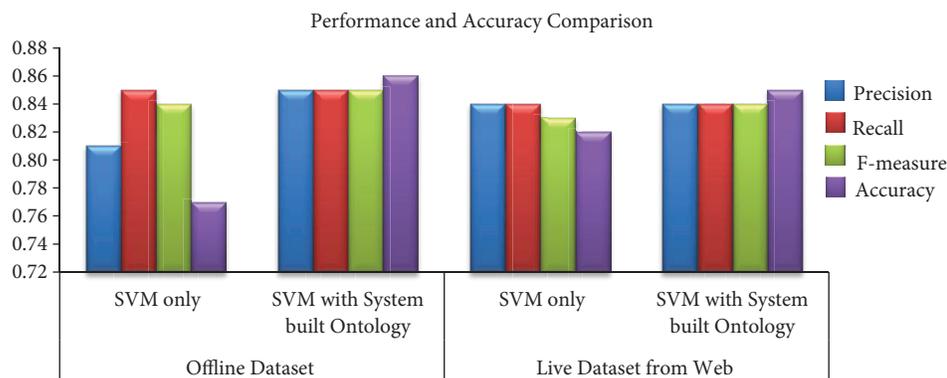


Figure 3. Comparison of performance and accuracy of two classifiers.

High accuracy of the system is also evident from the convex ROC under AUC analysis of two classifiers as shown in Figure 4 (Use-Case1) and Figure 5 (Use-Case2). ROC curves of both use-cases highlights that the SVM classifier with system-built ontology provides more accuracy vis à vis the simple SVM-based classifier. The ROC for the SVM only under Use-Case1 is quite uneven while the ROC for the SVM with system-built ontology provides a smooth convex curve, which is very promising. Both ROCs for Use-Case2 are smooth, but the curve for SVM with system-built ontology is more convex than the other.

Table 4. MCC metrics.

Average MCC value		
	SVM classifier	SVM classifier with system-built ontology
Use-Case1	0.33	0.62
Use-Case2	0.63	0.69

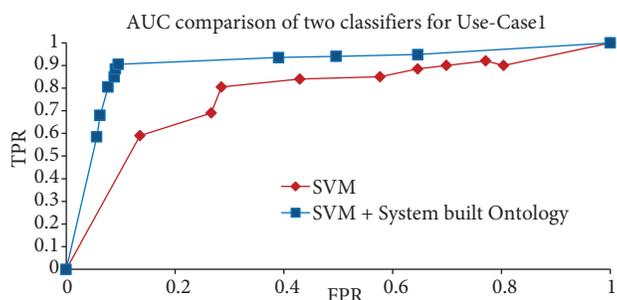


Figure 4. AUC comparison of two classifiers for Use-Case1.

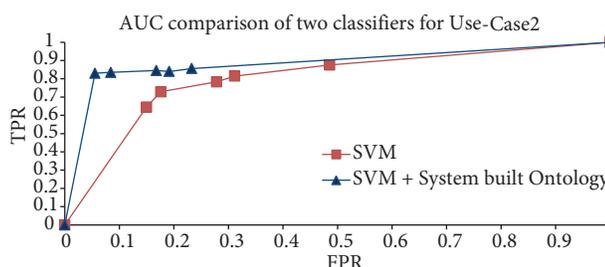


Figure 5. AUC comparison of two classifiers for Use-Case2.

Overall, the proposed self-governed ontology-based classification system provides us very favorable results, particularly in scenarios where the user feeds a natural language query to the system instead of a single word query.

6. Comparison with other methodologies

Our framework is similar to techniques such as those suggested by Luong et al. [12] in the manner that the ontology is learned and populated. Both frameworks are similar in terms of many points like using ‘focused crawling’ for retrieving documents from the web and automated ontology learning processes with the SVM as a classification method. However, there are a few critical differences in both the frameworks. The first and foremost difference is the use of ‘seed ontology’. While [12] uses a seed ontology, our framework starts from scratch and builds the seed ontology from the user queries. This highlights a step further in building a self-governed ontology-based classification system. The second and more critical difference between these two frameworks is the manner in which queries are generated and processed. While [12] implements a total automated query generation from the seed ontology to populate/enrich it further, our framework works in accordance with the user query, builds the seed ontology, and further enriches it automatically. Thus, our framework is more towards serving real-world users in the forefront, handling their queries at run time and providing them with relevant results as compared to [12], which serves for background enrichment of ontologies that will help users in the future while finding more appropriate results for their queries. The third difference between the two frameworks is their focus area of learning. While [12] is focused towards the biological domain

area, our framework is independent of any specific domain area and purely focused towards user query, picking up the domain at run time. Prefixing of the domain also requires [12] to generate seed ontology first and then proceed towards the enrichment process. Our framework does not have such dependencies, which results in making it a ‘self-governed’ learning system.

The experimental analysis of both systems also highlights quite a few contrasts. While [12] used a threshold of 0.6 during the experimental set-up, our framework uses a threshold of 1. This implies that [12] considered the top 60% of results, while we have considered the complete 100% results during performance and accuracy parameter calculations. If we use some threshold value, it will definitely help us improve our precision and accuracy more than what we have achieved under the current set-up. The overall F-measure and precision parameters of [12] (precision 88%, F-measure 85%) are slightly better than ours (precision 85%, F-measure 84%). This clearly shows a scope of improvement for our framework, which we shall try to achieve by incorporating a suitable threshold parameter in our experimental set-up.

7. Conclusion and future work

In this PoC, we propose a self-governed ontology-based approach to classify documents purely in the relevant context of the user query. The whole system gets triggered when a user throws a query in plain English language. Two branches of the system start working in parallel, one being the collection of relevant documents through a focused web crawler and the second being the build-up of seed ontology. Subsequent to data preprocessing and feature selection, again two processes start working in parallel, one being populating the domain ontology on top of seed ontology with the help of the ontology manager, and the second being the training of the SVM-based classifier. Towards the end, we compare the two setups for this system and conclude that the ontology-based classifier shows promising results and performs better than the simple SVM-based classifier. The whole system is implemented from scratch with no manually prepared seed ontology being used, which is quite encouraging. We have also compared our framework with similar available frameworks and found better usefulness of our framework in terms of self-governed learning system.

There are certain areas in the system design that point to our future work. First and foremost is the query manager, which handles the query string fed by the user and converts it into seed ontology by forming the RDF relation graph. Fine-tuning of this stage will provide better results. The OWL may also be explored during automatic ontology population instead of RDF. The second area of improvement is feature extraction and selection. We need to explore more mature algorithms during this stage. It will be quite interesting to perform the experiment using semantic kernels with the SVM instead of linear ones [24–26].

References

- [1] Christopher CS, Tylman J. Enterprise information portals. *Electron Libr* 1998; 18: 354–362.
- [2] Sheth A. Computing for human experience: semantics-empowered sensors, services, and social computing on the ubiquitous Web. *IEEE Internet Comput* 2010; 14: 88–91.
- [3] Sebastiani F. Machine learning in automated text categorization. *Comput Surv* 2002; 34: 1–47.
- [4] Gupta V, Lehal G. A survey of text mining techniques and applications. *Journal of Emerging Technologies in Web Intelligence* 2009; 1: 60–76.
- [5] Navigli R, Faralli S, Soroa A, de Lacalle OL, Agirre E. Two birds with one stone: learning semantic models for text categorization and word sense disambiguation. In: *20th ACM Conference on Information and Knowledge Management*; 24–28 October 2011; Glasgow, UK. New York, NY, USA: ACM. pp. 2317–2320.

- [6] Bizer C, Heath T, Berners-Lee T. Linked data—the story so far. *Int J Semant Web Inf* 2009; 5: 1–22.
- [7] Heath T, Bizer C. *Linked Data: Evolving the Web into a Global Data Space*. San Rafael, CA USA: Morgan & Claypool Publishers, 2011.
- [8] Buitelaar P, Cimiano P, Magnini B. *Ontology Learning from Text: Methods, Evaluation and Applications*. Philadelphia, PA, USA: IOS Press, 2005.
- [9] Maedche A, Staab S. Ontology learning for the semantic web. *IEEE Intell Syst* 2001; 16: 72–79.
- [10] Navigli R, Velardi P, Gangemi A. Ontology learning and its application to automated terminology translation. *IEEE Intell Syst* 2003; 18: 22–31.
- [11] Wei GY, Wu GX, Gu YY, Ling Y. An ontology based approach for chinese web texts classification. *Inform Technol J* 2008; 7: 796–801.
- [12] Luong HP, Gauch S, Wang Q. Ontology learning through focused crawling and information extraction. In: *International Conference on Knowledge and Systems Engineering*; 13–17 October 2009; Hanoi, Vietnam. New York, NY, USA: IEEE. pp. 106–112.
- [13] Brank J, Mladenic D, Grobelnik M. Large-scale hierarchical text classification using SVM and coding matrices. In: *Large-Scale Hierarchical Classification Workshop of the ECIR 2010*; 28–31 March 2010; Milton Keynes, UK.
- [14] Speretta M, Gauch S. Using text mining to enrich the vocabulary of domain ontologies. In: *IEEE International Conference on Web Intelligence and Intelligent Agent Technology*; 9–12 December 2008; Sydney, Australia. New York, NY, USA: IEEE. pp. 549–552.
- [15] Salton G, Buckley C. Term-weighting approaches in automatic text retrieval. *Inform Process Manag* 1998; 24: 513–523.
- [16] Abe S. *Support vector machines for pattern classification*. New York, NY, USA: Springer-Verlag, 2005.
- [17] RDF Working Group. *Resource Description Framework (RDF)*. W3C–Semantic Web, 2004.
- [18] Chang C, Lin C. *LIBSVM: A Library for Support Vector Machines*. Software, 2001.
- [19] Kohavi R. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: *International Joint Conference on Artificial Intelligence*; 1995; Montreal, Canada. pp. 1137–1145.
- [20] Matthews BW. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta* 1975; 405: 442–451.
- [21] Lewis D. *The Reuters-21578 Text Categorization Test Collection*, 1997.
- [22] Lang, K. NewsWeeder: Learning to filter netnews. In: *12th International Conference on Machine Learning*; 9–12 July 1995; Lake Tahoe, CA, USA. Washington, DC, USA: IMLS. pp. 331–339.
- [23] Craven M, DiPasquo D, Freitag D, McCallum A, Mitchell T, Nigam K, Slattery S. Learning to construct knowledge bases from the World Wide Web. *Artif Intell* 2000; 118: 69–114.
- [24] Siolas G, d’Alche Buc F. Support vector machines based on semantic kernel for text categorization. In: *International Joint Conference on Neural Networks*; 27 July 2000; Como, Italy. New York, NY, USA: IEEE. pp. 205–209.
- [25] Cristianini N, Shawe-Taylor J, Lodhi H. Latent semantic kernels. *J Intell Inf Syst* 2002; 18: 127–152.
- [26] Wang P, Domeniconi C. Building semantic kernels for text classification using Wikipedia. In: *14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*; 24–27 August 2008; Las Vegas, NV, USA. New York, NY, USA: ACM Press. pp. 713–721.