

Automatic Online News Gathering and Classification

Mr.Sachin Vinchurkar¹ and Prof.S.D.Bandari²

¹ Student of B.E. Information Technology
Dept. of Information Technology
DACOE,Karad,India

² Faculty of B.E. Information Technology
Dept.of Information Technology
DACOE,Karad,India

Acknowledgment

I have taken efforts in this project. However, it would not have been possible without the kind support and help of many individuals and organizations. I would like to extend my sincere thanks to all of them. I am highly indebted to my Guide Prof.S.D.Bandari for her guidance and constant supervision as well as for providing necessary information regarding the project & also for their support in completing the project. I would like to express my gratitude towards my parents & member of my institute for their kind co-operation and encouragement which help me in completion of this project. My thanks and appreciations also go to my colleague in developing the project and people who have willingly helped me out with their abilities.

Abstract: *In our proposed Categorizer system, we have experimented an automated approach to classify online news using the Naïve Bayes Algorithm. Naïve Bayes has been shown to deliver good classification results when sample training documents are given. In the proposed system, we have applied Naïve Bayes to the classification of online news.*

In the classification, users can define their personalized categories using a few keywords. By constructing search queries using these keywords, Categorizer obtains both positive and negative training documents required for the construction of personalized classifiers.

These keywords are known as the category profile for the newly created category. There is no restriction on the number of personalized categories for each user. To build the classifier for a personalized news category, a number of training documents have to be obtained. Instead of getting the user to perform the time-consuming task of selecting and uploading the training documents, we construct a query to the Yahoo News Search Engine using the category profile, i.e. the keywords. The training documents are then selected from the most highly ranked resultant news articles from Yahoo News.

Keywords: *Support Vector Machine, Naïve Bayes Algorithm, Yahoo news search engine, News classification, News monitoring, Feature selection, Syndromic surveillance*

Introduction

Classification of online news, in the past, has often been done manually. In our proposed Categorize system, we have experimented an automated approach to classify online news using the Naïve Byes. Naive Byes has been shown to deliver good classification results when ample training documents are given. In our research, we have applied Naïve Byes to the personalized classification of online news. In personalized classification, users can define their personalized categories using a few keywords. By constructing search queries using these keywords, Categorize obtains both positive and negative training documents required for the construction of personalized classifiers. In this paper, we describe the preliminary version of Categorize and present its system architecture.

Text classification is the process of assigning text documents to one or more predefined categories. This allows users to find desired information faster by searching only the relevant categories and not the entire information space. The importance of text classification is even more apparent when the information space is huge such as the World Wide Web. Examples of web classification systems include Yahoo! directory and Google web directory. However, such classification services are carried out by human experts, and they do not scale up well with the growth rate of web pages on the Internet.

To automate the classification process, machine learning methods have been introduced. In a text classification method based on machine learning, classifiers are built (trained) with a set of training documents. The trained classifiers can, therefore, assign documents to their suitable categories.

Before we introduce our scheme, this section first reviews several categories of existing solutions and their relationships to our work.

An author has described the Classification of online news, in the past, has often been done manually. SVM Is used for classification results when training documents are given. In our research, we have applied SVM to the personalized classification of online news. By constructing search queries using these keywords, Categorizer obtains both positive and negative training documents required for the construction of personalized classifiers. Text classification is a well-studied problem. Several methods have been proposed and many of them can be directly applied to news classification as long as there exist a good set of training documents for each predefined category. [1]

On the other hand, personalized classification is a form of personalization and there are several existing ways to support personalization. In the collaborative filtering approach, each user is associated with a user profile. When the user profiles of two users are similar, news articles that are interested in by one of them will be automatically recommended to the other. In another personalization approach known as content filtering, one or more sets of features each representing a different interest domain (personalized category) of a user are derived. News articles are then recommended based on the semantic similarity with each set of features. In this approach, the interesting domain of a user is very much independent of that of another user.[2]The Framework is unique and none of the techniques adopted in this study have been

previously used in the context of syndromic surveillance on infectious diseases. In the recent classification experiment, we compared the performance of different feature subset on different machine algorithm the result shown that the combined features subsets including Bags of Words Noun Phrases and Named entity feature out performed the Bags of Words feature subsets. [3]

These days, important public health related news is increasingly available on the World Wide Web in electronic form and has been shown to be a useful data source for syndromic surveillance. However, the volume of news is very large and there is a question on how to most effectively use this kind of information for syndromic surveillance to accurately detect the signals indicative of disease outbreaks. Syndromic surveillance systems that include a classification component can facilitate follow-up analysis and outbreak detection. However, to our knowledge, there is currently no automatic online news monitoring and classification system specifically designed for specific infectious diseases. It is also not clear what kind of document representation approach and machine learning algorithm perform best on online news classification for syndromic surveillance.[4]

This study is aimed at designing and examining automatic online news monitoring and classification methods for syndromic surveillance and global situational awareness. It provides an overview of literature concerning syndromic surveillance, news-based syndromic surveillance systems for infectious diseases, online data acquisition, text document representation, and feature selection used in text classification. We describe our research questions. We outline our architecture for automatic news monitoring and classification, after which we present our experiments and the results on foot-and-mouth disease (FMD) related online news. Finally, we describe our conclusions and future directions.[5]

Collecting data is a critical early step when developing a syndromic surveillance system. Data sources used in syndromic surveillance systems are expected to provide timely pre-diagnosis health indicators and are typically electronically stored and transmitted. Different types of data used for syndromic surveillance typically include: emergency department (ED) visit chief complaints, ambulatory visit records, hospital admissions, over-the-counter (OTC) drug sales from pharmacy stores, triage nurse calls, 911 calls, work or school absenteeism data, veterinary health records, laboratory test orders, and health department requests for influenza testing.[6]

Materials and Methods

The Naive Bayesian classifier is based on Bayes' theorem with independence assumptions between predictors.

We can see an example for this Naïve Bayes as follow:

Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. For example, a fruit may be considered to be an apple if it is red, round, and about 3 inches in diameter.

Figure .4. Example for Naive Bayes classifier

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

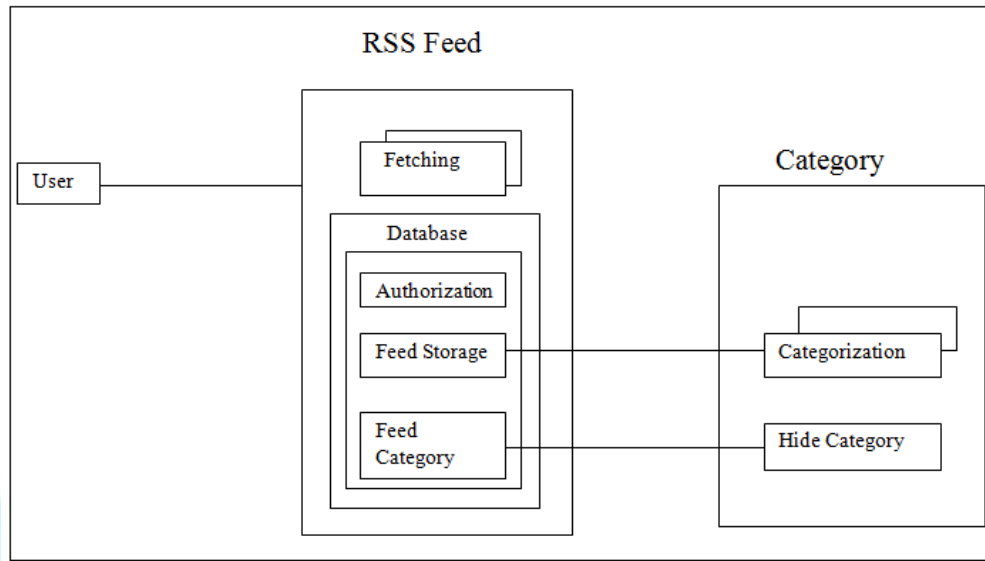
Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

1.3 System Architecture

System Architecture is a response to the conceptual and practical difficulties of the description and design of complex systems. Our architecture consists of main two modules RSS Feed and Category. A user is interacting with the system. RSS Feed block consists of a Fetch-ing process which can be executed multiple times. Authorization, feed storage, and feed category information are stored in the database. Category module consists of categorization and hides the category. Categorization is done multiple times as the feeds arrive. The line between Feed Storage and categorization indicates that they communicate between each other.

Figure .3. System Architecture

**Features:**

- 1.The unique feature of Categorizer is that it allows users to create and maintain their personalized categories.
- 2.Users can create their personalized news category by specifying few keywords associated with it. These keywords are known as the category profile for the newly created category.
3. There is no restriction on the number of personalized categories for each user.
4. To build the classifier for a personalized news category, a number of training documents have to be obtained.
5. Instead of getting the user to perform the time-consuming task of selecting and uploading the training documents, we construct a query to the Yahoo News Search Engine using the category profile, i.e., the keywords.
- 6.The training documents are then selected from the most highly ranked resultant news articles from Yahoo news.

Drawbacks of existing system:

- 1.The current updates are fetched only from the URLs listed under each category or URL given by users and not from all websites.
2. It is costly for the user..

Benefits of Proposed system:

- 1.To give current updates without time consumption is a key objective of this project.
- 2.The user can view the feeds of his interested categories.
- 3.There is no restriction on the number of personalized categories for each user.

Mathematical Evaluations

An original COCOMO is actually a collection of three models: a basic model that can be applied when little about the project is known, an intermediate model that is applied when the design is complete. All three take the same form,

Formula: $E_i = a * (KLOC)^b$

$$E_i = 2.4 * (2.736)^{1.05}$$

$$E_i = 6.9026$$

$$E = E_i * EAF$$

$$E = 6.9026 * 1.00 * 1.08 * 1.00 * 1.00 * 1.00 * 1.00 * 0.91$$

$$E = 6.7838 \text{ pm}$$

Where E_i is the initial efforts, EAF is the factor, E is effort in person months, S is size measured in thousands of delivered source instructions (KDSI)

$$D = a * E^b$$

$$D = 2.5 * (6.7838)^{0.38}$$

$$D = 5.1748 \text{ months} \quad \text{Where, D is the duration in months.}$$

Hence, the project requires 7 persons per months to complete the software.

Result and Analysis

For Experimental results, we are going to calculate the values for Precision and FI-Score. We are taken text features like Sports, us, Sci/Tech, bus category, word, rent, and health. We can see all values are in following tables.

Table 1:Text feature,Precision,recall and FI-Score

Text Feature	Precision	Recall	F1-Score
Sport	0.75	0.85	0.8
US	0.59	0.62	0.6
Sci/Tech	0.65	0.55	0.6
Bus Category	0.61	0.6	0.6
World	0.77	0.75	0.77
Rent	0.6	0.62	0.61
Health	0.6	0.5	0.55

Conclusion And Future Scope

We have designed and implemented the news classification system based on the Naïve Bayes classification method. The system is capable of both general Classification and personalized classification. Firstly, we enhanced the Categorizer with a complete set of general categories due to the unavailability of generic extraction software for extracting the desired news text from HTML web pages. We have studied Naïve Byes algorithm which is used to text classification. This application is useful for everyone. It reduces the user time and provides required information in minimum time and effort. This application contains less complexity.

We can develop news gathering and classification system. In this system we can implement personalized categorization. In future our project will extract videos from internet such as movies, video etc. also user can be get live updated news from the internet. User will be added his own interested news and share that news with another user. User will be personalized his own category as per his demand and provide more security for each user. And also we will try to use regional languages which are user friendly.

References

1. Chee-Hong Chan, et.al, Automated Online News Classification with Personalization, INTERNATIONAL CONFERENCE OF ASIAN DIGITAL LIBRARY (ICADL), VOL.2, PP.320-329, 2001.
2. Yulei Zhang, et.al, Automatic Online News Monitoring and Classification for syndromic Surveillance, VOL.3, PP.1-10, 2009.

3. S. Dumais and H. Chen. Hierarchical classification of Web content. 23rd ACM INTERNATIONAL CONFERENCE ON RESEARCH AND DEVELOPMENT IN INFORMATION RETRIEVAL, PP.256-263, 2000.
4. S. Dumais, et al, Inductive learning algorithms and representations for text categorization. 7th INTERNATIONAL CONFERENCE ON INFORMATION AND KNOWLEDGE MANAGEMENT, PP.148-155, 1998.
5. Susan Dumais and Hao Chen, HIERARCHICAL CLASSIFICATION OF WEB CONTENT, PP.1-9, 2000.
6. Susan Dumais, et al, Inductive Learning Algorithms and Representations for Text Categorization, PP.1-8, 2001.
7. Mark A. Hall, Correlation based feature selection for discrete and numeric class machine learning, University of Waikato, Hamilton, New Zealand, PP.1-8, 1997.
8. M. Dash and H. Liu, Feature Selection for Classification, Intelligent Data Analysis, PP. 131-156, 1997.
9. Thorsten Joachims, Text categorization with support vector machines: Learning with many relevant features, PP.1-7, 1998
10. Channel News Asia, <http://www.channelnewsasia.com>.
11. <http://www.howtogeek.com>
12. <http://www.ampercent.com>
13. Yahoo! News, <http://news.yahoo.com>.