

# A Semi-Automatic Emerging Technology Trend Classifier Using SCOPUS and PATSTAT

Seonho Kim, Woondong Yeo, Byong-Youl Coh, Waqas Rasheed<sup>†</sup>, Jaewoo Kang<sup>†</sup>

*Technology Information Analysis Lab. Department of Information Analysis, Korea Institute of Science and Technology Information, 66, Hoegi-ro, Dongdaemun-gu, Seoul 130-741, Korea*  
{haebang, wdyeo, cohby}(at)kisti.re.kr

<sup>†</sup> *Data Mining and Information System Lab. College of Information and Communication, Korea University, 5-ga, Anam-dong, Seongbuk-gu, Seoul 136-713, Korea*  
{waqas, kangj}(at)korea.ac.kr

**Abstract:** *Identifying Emerging Technology Trends is crucial for decision makers of nations and organizations in order to use limited resources, such as time, money, etc., efficiently. Many researchers have proposed emerging trend detection systems based on a popularity analysis of the document, but this still needs to be improved.*

*In this paper, an emerging trend detection classifier is proposed which uses both academic and industrial data, SCOPUS [1] and PATSTAT [2]. Unlike most previous research, our emerging technology trend classifier utilizes supervised, semi-automatic, machine learning techniques to improve the precision of the results. In addition, the citation information from among the SCOPUS data is analyzed to identify the early signals of emerging technology trends.*

**Keywords:** *SCOPUS, PATSTAT, Emerging trend detection, Machine learning, Artificial neural network*

## I. INTRODUCTION

Research on emerging trend detection has been conducted by many researchers and organizations, because identifying emerging trends and predicting the near future can improve efficient resource distribution and effective policy establishment. Traditional research on detecting emerging trends is based on massive data analysis, literature review, and brainstorming of intellectuals [3]. However, as the amount of information to analyze increases and computing technology develops, recent trends of this research are focused on automating the processing by utilizing text mining and data analysis techniques [4, 5]. In this study, we propose a semi-automatic machine learning technology for recognizing emerging technology trends from academic and industrial text data.

## II. METHODS

### A. Previous Research

Traditional research on emerging trend detection in text data is conducted by two major methods, fully automatic detection and semi-automatic detection. The fully automatic methods, such as TimeMines [6], VUDM [7], and on-line event detection by Allan [8], have employed unsupervised learning methods, in which learning is achieved in a fully automatic manner

by detecting topics describing technologies to analyze their trends. These approaches tend to show good performance in recall, which represents the portion of emerging trends found over the total number of emerging trends in the system, but only medium performance in precision, which expresses the rate of correctly determined emerging trends over the total number of trends found. The reason for the good recall of the method is that it uses bottom-up methods to test every possible combination of noun phrases as topics describing emerging trends. After finding topics, the trends of each topic are analyzed, and emerging trends are selected from among the topics based on certain criteria. However, the results of this approach can be noisy for some decision makers because, in most cases, the identification of several correct targets is more meaningful than identifying as many correct targets as possible without missing.

For this reason, in order to increase precision rather than the recall, semi-automatic methods which use supervised learning are attracting attention. These approaches, such as CIMEL [9], PatentMiner [10], and HDDI [11], use experts to provide guidance to the machine along with the training data. For example, the training data of the machine learning is tagged by the experts to identify whether or not the data is about an emerging trend, or the experts interact with the software during learning. Therefore, the machine can learn more correct patterns of emerging trends and classify new trends more precisely.

### B. Semi-automatic Supervised Machine Learning

Our proposed method is a semi-automatic system in which experts provide guidance to an artificial neural network machine learning system. The guidance from the experts in our system is the classification tags attached to the training data, so that the machine is directed as to whether or not the information currently being learned is about emerging technology. In a machine learning system, selecting the correct features for learning is most critical. HDDI [11] used the frequencies of concepts within a fixed length of periods and the count of concepts in semantically defined regions, which can be obtained automatically. Oh [12] also provided the results of a comparison experiment for testing the strength of impact of the most famous features of trends generally used for emerging trend

detection, such as change, persistency, stability, and volume. According to Oh's result, the strongest impact feature of a trend is expressed by the change, stability, volume, and persistency in descending order. However, because every emerging trend detection systems have different purposes, available source data, and training methods, the features to be learned need to be selected and tested themselves. Our system uses both academic data, SCOPUS [1], and industrial data, PATSTAT [2]. Thus, in addition to the usual term-frequency based trend changes in academic data, the latency between the weak academic signal and the industrial legal activity (patent application) could be trained. Also, our system uses a citation half-life feature, which is described in the next section, to train the aging factor of technology. We assume that the age of the emerging technology is relatively young and young technology tends to have a later citation half-life than mature technology.

This study consists of three stages, 1) emerging trend detection, 2) training data generation, and 3) learning and experimentation. Figure 1 shows the brief structure of the proposed system.

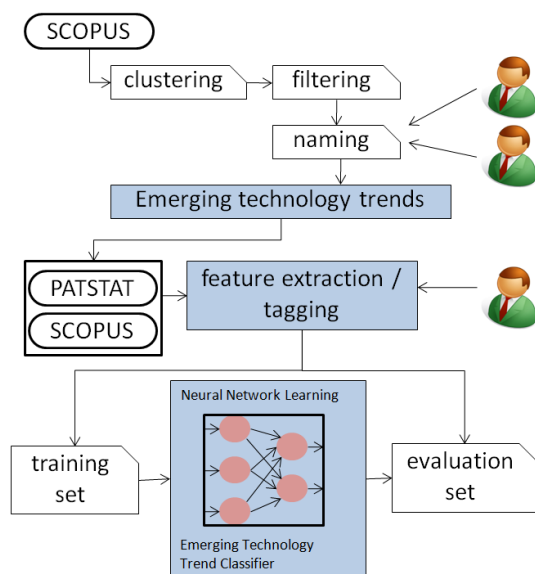


Figure 1. Emerging technology trend classifier

At stage 1, a set of emerging technology trends are identified from SCOPUS. The process contains clustering, filtering, and naming. The field experts intervene during the process. After that, in the second stage, the training data is generated by extracting learning features from SCOPUS and PATSTAT. Experts also join in this process to add their guidance to the learning data.

### C. Emerging Technology Trends

Detecting emerging trend from SCOPUS is conducted to generate tagged feature vectors of training data for neural network learning. Only 1% of the highly cited papers are used. The papers are clustered using the scientometrics techniques, which use the citation

information, classification code, core paper identification, and shared term frequency. As a result of clustering, 512 clusters are identified. Then, about 50 field experts are selected for the naming of the clusters and voting on the prospects of the technologies as expressed by their names. Clusters that do not have any consistent topic, based on the knowledge of the field experts, are filtered out, and final 87 clusters are selected as emerging technology trends. From this result, 60 technologies are selected for positive training data. In addition, 26 matured, non-emerging, technologies, such as compact discs, fountain pens, optical character recognition, and bluetooth technology, are prepared as negative training data.

### D. Feature Extraction and Tagging

An artificial neural network is trained with positive training data, which consist of features from emerging technology trends, and negative training data, which consist of features from matured, non-emerging, technology trends.

The training features are extracted from SCOPUS and PATSTAT data from 1996 to 2010. Two types of feature, linguistic and statistical, are used. For the linguistic features, the trend "change" for 2 different periods from SCOPUS and PATSTAT are used. For statistical features, the citation "half-life" in SCOPUS is used. The change value represents the slope of the change between the start and end point of the trend curve. A positive change value means that the trend is emerging, while a negative change value means that the trend is submerging [12].

$$\text{diff}(f) = LR_f(T_{end}) - LR_f(T_{start}) \quad (1)$$

$$\text{Change}(f) = \frac{\text{diff}(f)}{\sqrt{\text{diff}(f)^2 + (T_{end} - T_{start})^2}} \quad (2)$$

The change feature is obtained by calculating the linear regression, which is represented by equations (1) and (2). In these equations,  $f$  is a trend function expressed by an independent value, time, and a dependent value, the populations in SCOPUS and PATSTAT.  $LR_f$  is the linear regression function about  $f$ .  $T_{end}$  expresses the population at the end point of the period of interest, and  $T_{start}$  is the population at the start point of the period. For the machine learning part of this study, four change features, 1) the change in the first 11 years in SCOPUS, 2) the change in the last 4 years in SCOPUS, 3) the change in the first 11 years in PATSTAT, and 4) the change in the last 4 years in PATSTAT, are extracted from one technology trend curve.

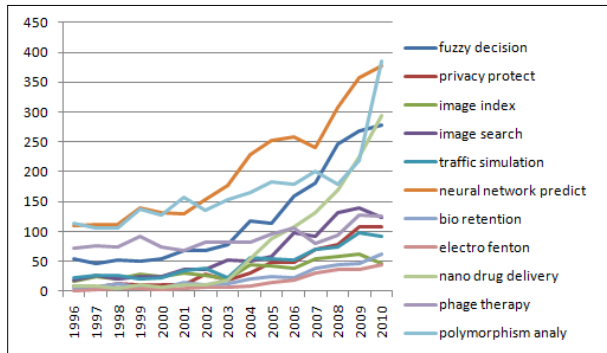


Figure 2. Emerging technologies' population in SCOPUS

The reason for dividing the period into two sections, 11 years and 4 years, is based on our basic interpretation about the data. Figure 2 and figure 3 show some examples of trend curves of emerging technologies in SCOPUS and PATSTAT, respectively, and both charts show that the trend curves changed noticeably in the last 4~5 years.

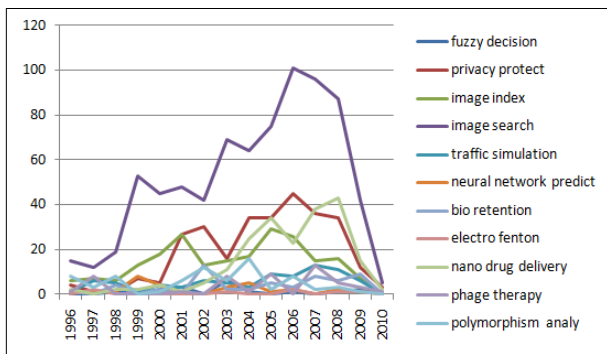


Figure 3. Emerging technologies' population in PATSTAT

According to the figure 2, the population of emerging technology trends had increased steadily in the first 11 years and its changes escalated during the last 4 years in SCOPUS. In figure 3 the speed of change dropped rapidly during the last 4 years in PATSTAT.

The other type of feature selected for this study is the citation half-life, which is the number of years that have passed after half of the total number of citations concerning a particular technology. The citation half-life represents the center of the total group of citations. Figure 4 illustrates the concept of citation half-lifetime.

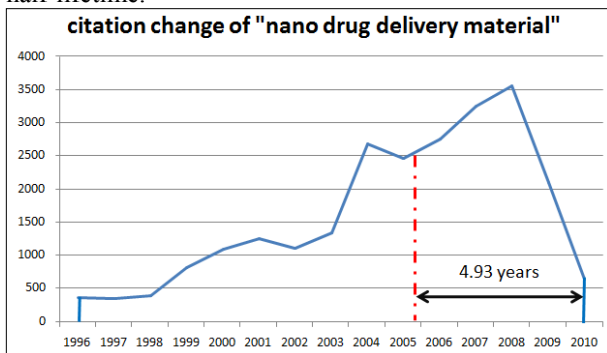


Figure 4. the concept of citation half-life

The number of citation on "nano drug delivery materials" increased steadily until its peak in 2008, and then it decreased rapidly until the present year, 2010. The total number of citations on the technology during the observation period in SCOPUS is 24,225 and the half of the total number of citations, 12,112.5, is met around 2005. Therefore, the citation half-life of the "nano drug delivery material" technology is 4.93 years.

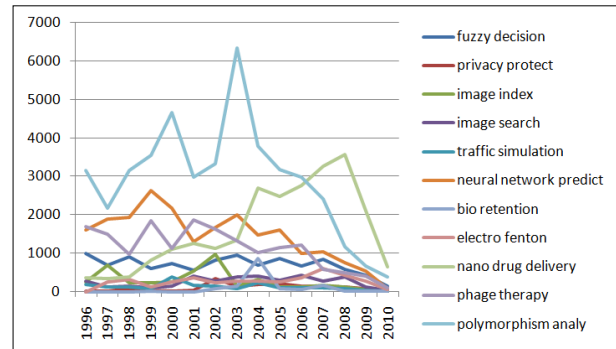


Figure 5. Emerging technologies' citation number in SCOPUS

Every technology has a distinct trend in citation. Figure 5 shows the citation trends of some example emerging technologies. The trend curves tend to bias to the left or right based on the age of the technology. Thus, we assumed that an emerging technology is a relatively new technology, and the citation half-life value represents the age of technology, and that learning this value may help detect emerging technologies.

### E. Learning

The goal of machine learning in this study is to educate the neural network with the features of emerging technologies and non-emerging technologies, and evaluate the performance of the trained neural network as an emerging technology trend classifier.

During learning, a total 86 feature vectors, 60 from positive data and 26 from negative data, and tagged as to whether or not it is emerging technology, are entered into the neural network for training.

In the evaluation session, the precision of the classification performance is measured. In order to evaluate the precision of the learned neural network, 64 new feature vectors, which have not been used for training, are entered. The evaluation data is not tagged but are already classified by field experts and compared with the decisions made by the neural network.

### F. Evaluation and Results

The final status of a neural network, the learned status, is different in every trial, even if it is trained with identical training data, because the neural network assigns random initial weights to the links among the neurons in the hidden layers. Therefore, we repeated the training and evaluations 10 thousand times, with the same evaluation data, and determined the average precision and a standard deviation.

The average precision of classification of the neural

network is 69.92 and the standard deviation is 21.93. Regarding precision only, this result is better than the normal precision rate, which is about 45, for normal unsupervised learning for emerging trend technology.

The precision of learning can be improved by iterative tuning of the training data by filtering out noisy data and adding correct data, and the deviation can be decreased by enlarging the size of training data. However, the tuning process is not included in this study.

### III. CONCLUSIONS AND FUTURE WORK

An emerging technology trend identification system is proposed, which is based on semi-automatic supervised machine learning technology, and its performance is evaluated. Our method learns from both academic data, SCOPUS, and industrial data, PATSTAT. The change in trend and citation half-life is selected as the learning feature. By using a supervised learning method, our system places the focus on enhancing the precision of classification rather than recall. This system can be applied for fast classification of emerging technology, such as a real time online technology evaluator or competitive intelligence system.

For future work, we will tune the training data to make the learning process more efficient. Also, a larger volume of training data and evaluation data may improve the precision of the classifier. Furthermore, we plan to apply the standard test set for emerging trend detection, such as the Topic Detection and Tracking (TDT) corpora [13], to make our system applicable as a general purpose topic detection system. Finally, we will test more various trend features, such as changes and statistics about search queries.

### ACKNOWLEDGEMENTS

Thanks go to Korea Ministry of Education, Science And Technology for support of grants K-11-L05-C03-S02 for this research.

### REFERENCES

- Elsevier. *SCOPUS*. 2011; Available from: <http://info.scopus.com/>.
- PATSTAT. *EPO Worldwide Patent Statistical Database*. 2011 [cited 2011; Available from: <http://www.epo.org/searching/subscription/raw/product-14-24.html>].
- Shibata, N., Y. Kajikawa, and K. Matsushima. *Detecting emerging research fronts based on topological measures in citation networks of scientific publications*. *Technovation*, November 2008. **28**(11): p. 758-775.
- <http://www.mendeley.com/research/detecting-emerging-research-fronts-based-on-topological-measures-in-citation-networks-of-scientific-publications/>
- Smalheiser, N.R., *Predicting emerging technologies with the aid of text-based data mining: the micro approach*. *Technovation*, October 2001. **21**(10): p. 689-693.
- <http://www.sciencedirect.com/science/article/pii/S0166497201000487>
- Kontostathis, A., et al., *A Survey of Emerging Trend Detection in Textual Data Mining*. 2003.
- <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.3.6551&rep=rep1&type=pdf>
- Swan, R. and D. Jensen. *Timemines: Constructing timelines with statistical models of word usage*. in *the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2000. Boston, MA, USA: Association for Computing Machinery.
- [http://www.cs.cmu.edu/~dunja/KDDpapers/Swan\\_TM.pdf](http://www.cs.cmu.edu/~dunja/KDDpapers/Swan_TM.pdf)
- Kim, S., *Visualizing Users, User Communities, and Usage Trends in Complex Information Systems Using Implicit Rating Data*, Dissertation, Department of Computer Science, Spring 2008, Virginia Tech: Blacksburg, VA, USA.
- <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.105.841&rep=rep1&type=pdf>
- Allan, J., R. Papka, and V. Lavrenko. *On-line New Event Detection and Tracking*. in *the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1998. Melbourne, Australia: Association for Computing Machinery.
- <http://web.cs.dal.ca/~watters/courses2005/6403/p37-allan.pdf>
- Roy, S., D. Gevry, and W.M. Pottenger. *Methodologies for Trend Detection in Textual Data Mining*. in *the Textmine 2002 Workshop, Second SIAM International Conference on Data Mining*. April 2002. Arlington, VA.
- <http://dimacs.rutgers.edu/~billp/pubs/ETDMethodologies.pdf>
- Lent, B., R. Agrawal, and R. Srikant, *Discovering Trends in Text Databases*. 1997: AAAI Press.
- <http://www.aaai.org/Papers/KDD/1997/KDD97-046.pdf>
- Pottenger, W.M., Y.-B. Kim, and D.D. Meling, *HDDI: Hierarchical Distributed Dynamic Indexing*, in *Data Mining for Scientific and Engineering Applications*, R.L. Grossman, Editor. 2001, Springer-Verlag: New York, LLC.
- <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.18.9475>
- Oh, H.-S., et al., *Trend Properties and a Ranking Method for Automatic Trend Analysis*. *Journal of KIISE: Software and Applications*, March 2009. **36**(3): p. 236-243.
- <http://ir.kaist.ac.kr/anthology/2009.03-%EC%98%A4%ED%9D%A5%EC%84%A0.pdf>
- NIST. *TDT. Topic Detection and Tracking Evaluation*. 2011; Available from: <http://www.itl.nist.gov/iad/mig//tests/tdt/index.html>