# Discovering the Dynamics of Terms' Semantic Relatedness through Twitter

Nikola Milikic[1], Jelena Jovanovic[1], Milan Stankovic[2]

[1]University of Belgrade, Jove Ilica 154, 11000 Belgrade, Serbia
[2]STIH, Université Paris-Sorbonne, 28 rue Serpente, 75006 Paris, France

nikola.milikic@gmail.com, jeljov@gmail.com, milstan@gmail.com

**Abstract.** Determining the semantic relatedness (SR) of two terms has been an appealing topic in information retrieval for many years as such information is useful for various tasks ranging from tag recommendation, over search query refinement to suggesting new web resources for the user to discover. Most approaches consider the SR of terms as static over time, and disregard the eventual temporal changes as imperfections. However, detecting and tracing changes in SR of terms over time may help in understanding the nature of changes in public opinion, as well as the change in the usage of terms in common language and jargon. In this paper, we propose an approach that makes use of microposts data in order to establish a dynamic measure of SR of terms, i.e., a measure that accounts for the changes in SR over time. We propose different scenarios of use (in online advertising and organizational knowledge management) which demonstrate the applicability of our approach in real life situations. We also provide a demo application for visualizing the change in micropost-based SR of terms.

**Keywords:** Semantic relatedness, dynamic measure of semantic relatedness, microposts, Twitter

## 1 Introduction

Many research papers (such as in Wagner [1]) claim that Twitter and similar micro-blogging services have become a valuable source of knowledge, and have tried to extract this knowledge and use it for various purposes, such as the creation of dynamic domain models suitable for semantic analysis and annotation of real-time data [2], modeling of users' interests and finding experts [3], etc. However, from our point of view, the real-time nature of Twitter and Twitter-like services has not really been explored to its full extent, yet.

Most approaches exploit the mass of data that users generate on real-time services as their most valuable feature. We believe that there is a significant value in the fact that tweets (and microposts in general), posted frequently and massively, represent the moment in which they are created and the characteristics of that moment. Therefore, we have been exploring how these real-time services can support the detection of changes

in semantics of terms, by enabling one to observe the changes of a term's use over time. We focus particularly on the semantic relatedness (SR) of terms which is also subject to temporal changes.

The scope of meaning of a certain term is always defined in a social circle in which that meaning emerges, is agreed upon and accepted. Knowing that social systems are dynamic, it is difficult to neglect the natural changes (i.e., evolution) in socially agreed upon meaning of terms. If the meaning of a term is changing over time, so is the relatedness of that term to other terms whose dynamics in the given time period might be different. The most basic illustration of this is the term *totalitarian regime*. It is reasonably close to the terms identifying particular totalitarian governments and dictators of particular countries. However, this proximity should decrease if the totalitarian regime in a country is replaced by a democratic government – which happens more and more often in recent times.

Although tendencies in the public expressions can easily be detected through search query frequencies and trending topics on Twitter, the nature of tendencies and their mutual relationships are not directly evident from such observations. We could imagine having three trending topics on Twitter: *Egypt*, *revolution*, and *Britney Spears*. Although a human may grasp that it is more likely that the revolution is happening in Egypt and not that Britney Spears is leading a revolution, for a computer, this is far less obvious. The change of SR, however, could indicate the rationale for the raising public interest in a particular term. For instance, we could see that the recent popularity of the term *Egypt* might have been related to the temporary increase of SR of terms *Egypt* and *revolution*, and that it had nothing to do with the raise of popularity of the term *Britney Spears*. Spotting the change in SR of terms could thus help to give meaning to the observed trends in Web content, and enable machines to grasp this meaning and take advantage of it in many real life scenarios.

In this paper we present our initial research on using real-time services, in general and Twitter in particular, to detect the changes in SR of terms. We also explore the scenarios where reacting to those changes might be beneficial. In Section 2, we present the state of the art in research on SR of terms as well as in using Twitter to detect tendencies and make use of them. Section 3 introduces our measure of SR based on micropost data – Normalized Micropost Distance, whereas Section 4 gives some suggestions on how the relevancy of the change in SR of terms could be detected. We present our application for testing the proposed approach in Section 5, and consider the potential usage scenarios in Section 6. In Section 7, we give some interesting examples of changes in terms' SR which we have observed by using our application. Section 8 concludes the paper with propositions of future work that will help give maturity to our initial research.

## 2 State of The Art

The problem of determining semantic relatedness of terms has been studied for decades, in various contexts and using different approaches. Semantically related terms have been used to help users choose the right tags in collaborative filtering systems [4]; to

discover alternative search queries [5]; for query refinement [6]; to enhance expert finding results [7]; for ontology maintenance [8] [9], and in many other scenarios.

Different techniques and different sources have been used and combined to develop measures of semantic relatedness (MSRs). These measures could be split into three major categories: 1) net-based measures, 2) distributional measures and 3) Wikipedia-based measures [10]. In what follows we briefly examine each category of MSRs.

Net-based measures make use of semantic (e.g., hyponymy or meronymy) and/or lexical (e.g., synonyms) relationships within a network (graph) of concepts to determine semantic proximity between the concepts. For example, Burton-Jones et al. [11] exploit the hypernym graphs of WordNet[1]; Safar et al. [6] use Gallois lattice to provide recommendations based on domain ontologies, whereas Ziegler et al. [12] and Resnik [13] use the ODP taxonomy[2]. This category also includes measures that rely on the graph structure of concepts to determine semantic relatedness of those concepts. Shortest path is among the most common of such measures. It is often enhanced by taking into account the informational content of the nodes in the graph [14].

Distributional measures rely on the distributional properties of words in large text corpora. Such MSRs deduce semantic relatedness by leveraging co-occurrences of concepts. For example, the approach presented in Salton et al. [15] uses co-occurrence in text of research papers, pondered with a function derived from the tf-idf measure to establish a notion of word proximity. Co-occurrence in tags [4] and in search results [17] is also commonly used. In Strube et al. [18], the authors introduced Normalized Web Distance (NWD) as a generalization of Normalized Google Distance (NGD) MSR and investigated its performance with six different search engines. The evaluation (based on the correlation with human judgment) demonstrated the best performance of Exalead-based NWD measure, closely followed by Yahoo!, AltaVista, Ask and Google; only Live Search and Clusty showed significantly lower results.

As its name suggests, the third category of MSRs – Wikipedia-based measures – makes use of Wikipedia as the resource for computing semantic relatedness and often combines the features of the previous two MSR groups. For example, [18] relies on the graph of Wikipedia categories, whereas Waltinger et al. [10] rely on co-occurrence of words in the text of Wikipedia pages, combined with the information about the categories of pages in Wikipedia to compute semantic relatedness.

In Waltinger et al. [10], the authors report on a comparative analysis of a large number of MSRs (at least 4 algorithms from each major category of MSRs were included in the study, resulting in sixteen algorithms in total). The most important results could be summarized as follows: 1) small, hand-crafted and structured resources (e.g., WordNet) are inferior to large and semi-structured (i.e., Wikipedia) or even unstructured resources (i.e., plain text); 2) the distributional MSRs (especially measures like Latent Semantic Analysis) perform significantly better than the net-based measures and those using explicit categorical information; 3) MSRs that use the Web as a corpus were inferior to those operating on smaller but better controlled training corpora (e.g., Normalized Distance based on Wikipedia significantly outperformed NGD).

Most of the existing approaches do not take into account the dynamic nature of semantic relatedness between terms. An exception would be the work presented in

---

[1] http://wordnet.princeton.edu/

[2] http://www.dmoz.org/

Nagarajan et al. [22] where authors take the approach of identifying 'strong descriptors' of an event by querying Google Insights to get the terms the event's name was queried with the most (referred as 'seed keywords'). Afterwards, they query Twitter to get the tweets containing seed keywords and extract the strong descriptors from them. However, this approach does not measure SR between two specific terms, but rather identify terms relevant to the name of an event being examined. Other approaches even take the stability of their measure over time, to demonstrate the solidity of their approach [17].

On the other hand many approaches exist for extracting meaning from Twitter [20][1]. Some of them make extensive use of Twitter dynamics, like the approach for detecting events through peaks of word popularity [20]. Most related to our work is the approach presented in Song et al. [21] which relies on spatio-temporal characteristics of topics mined from Twitter data, for the calculation of semantic relatedness among topics. The temporal aspect of a topic is determined by the frequency of its occurrence in Twitter data streams over a given time period, whereas the spatial aspect refers to the regional distribution of messages mention the given topic over the same time period. Although this approach looks promising, its usefulness for measuring SR of topics has not been fully proved yet.

## 3  Normalized Micropost Distance

Inspired by the work of Cilibrasi et al. [16] on establishing Normalized Google Distance (NGD) as a MSR of terms based on Google search result, we propose a similar measure – Normalized Micropost Distance (NMD) – based on the results of searching the content (i.e., microposts) of real-time (Twitter-like) services. By leveraging micropost streams of real-time services, this measure should reflect the change in terms' SR more quickly than the standard web search results that are not updated in real-time. The basic assumption behind our approach is that Google's Search API results tend to be stable and based on content with a lower frequency of change, and as such would not be as good in indicating the changes in the SR of terms as could be search results that are based on real-time content..

NGD uses the frequencies of appearance of two terms in the Google index, as well as the frequency of their mutual appearance to quantify the extent to which the two terms are related. The basic assumption behind this measure is that terms that co-occur more frequently would be more related. Similarly, the proposed NMD measure can be calculated using the formula (1).

$$NMD(x,y)_t = \frac{\max\{\log f(x)_t, \log f(y)_t\} - \log f(x,y)_t}{\log M - \min\{\log f(x)_t, \log f(y)_t\}} \quad (1)$$

The formula allows one to calculate the NMD of two terms $x$ and $y$ for the time interval $t$. $f(x)_t$ and $f(y)_t$ represent the number of results returned for the term $x$ and $y$, respectively, within the time interval $t$, when searching the content (i.e., microposts) of a real-time, Twitter-like service. The terms in the formula ($x$ and $y$) may also be compound terms. Calculating the value of this formula for the same terms over different

time intervals is essential for determining the dynamics of their relationship, as we further explain in the following two sections.

## 4  Detecting the Significance of Change

The notion of NMD defined above is useful for measuring the difference in SR of two terms, but will not, by itself help to detect changes worthy of notice, and distinguish them from small and frequent variations. We suggest two complementary ways to perform this detection.

First, calculating the standard deviation of NMDs over a longer period of time would give a good ground to judging the significance of the identified changes. Standard deviation of NMDs can be calculated using the formula (2). The given formula represents the standard deviation of NMDs over a sample of N observations in which NMDs were calculated in time intervals $i$ that are of the same length.

$$\sigma(NMD(x, y)) = \sqrt{\frac{\sum_{i=1}^{N}(NMD(x, y)_i - avg(NMD(x, y)))^2}{N}} \tag{2}$$

Detection of a change in terms' SR (measured using NMD) that is greater than the standard deviation $\sigma$ could be an indicator of a significant change.

In addition to this indicator, one could observe the stability of change over several consecutive time instances to make sure that the change is not of a too short breath. However such a criterion may not be generally applicable and is specific to each use case, as even short changes might matter in some use cases, while in others only a change that spans several days would be significant.

## 5  Demo Application

In order to test the proposed approach of using micropost streams to calculate SR of terms, we have developed a simple web application that makes use of Twitter Search API[3] for computing NMD. The application, entitled Tweet Dynamics, currently in private beta, demonstrates how the NMD measure can be utilized, visualized and interpreted. Application is built in Java programming language using Tapestry Web Framework[4]. Javascript plotting library for jQuery named Flot[5] is used for plotting the result diagram.

The application's home page presents a user with a simple interface (Figure 1) which allows her to input the number of days and two keywords that NMDs should be calculated for. By clicking on the button 'Calculate', NMD calculation process is

---

[3] http://search.twitter.com/api/

[4] http://tapestry.apache.org/

[5] http://code.google.com/p/flot/

invoked. Application then queries the Twitter API to get all posts containing the first keyword, then posts containing the second keyword, and at the end to get all the posts containing both keyword. This process is repeated for the given number of days. With that data, NMDs are being calculated according to the formula (1).

The result of calculation is shown in a diagram (Figure 2) where each day is presented as a dot on the diagram line. One can easily perceive a trend of SR between two keywords during the past days.

Although, for the purpose of calculating standard deviation, our application keeps the computed values of NMD, the value of standard deviation is not shown on Figure 2 since we do not yet have a significant sample of values (e.g., dating from at least a month ago) and thus taking into account the currently available value of standard deviation would not be methodologically sound. Once a significant sample is present, the user would see a second line representing the standard deviation, so he/she could spot when the change in NMD becomes significant.
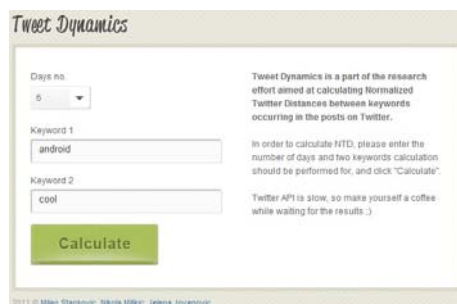


Figure 1 - Tweet Dynamics home page. User can enter two terms and a number of days to observe.



Figure 2 – The diagram illustrates the dynamics of SR of the terms *android* and *cool*, measured using NMD

Although some Web actors have access to the total history of tweets, most of interested parties have quite limited access to the Twitter Search API, which allows up to 1500 results per query. For terms of high frequency this can be a limiting factor since it makes it impossible to estimate their full frequency, and compare it with other high-frequency terms. A workaround that we use is to sample the tweets within short time intervals in which the number of tweets per terms is lower than the imposed limit. This however involves the risk of hitting the limit of 150 requests to the Twitter API per hour.

Another limitation of using Twitter's Search API is the restriction on the temporal range of tweets that can be returned as a search result. In particular, according to the API's documentation,[6] a post returned as a search result must not be 'too old', which in practice brings down to a number of six days, meaning that the oldest post returned as a result of a query is six days old. This restriction highly limits the ability to test our application on the microposts generated during a longer time span and detect trends in SR between keywords related to certain events or periods of year. If we had access to data spanning a longer period of time, we would have been able to test our results by

---

[6] http://apiwiki.twitter.com/w/page/22554756/Twitter-Search-API-Method:-search

comparing them with various indicators such as survey results, sales changes of a product etc.

# 6  Scenarios of Use

In this section we present two usage scenarios aiming to illustrate the potential benefits of the suggested dynamic MSR in real life settings. The first scenario assumes the usage of Twitter content stream for the calculation of NMD, whereas the second one relies on the micropost exchanged in a (internal) micro-blogging tool of an organization.

### Scenario 1: Adapting Online Advertising Campaigns to the Changes in Term Relatedness

Optimization of the keyword choice for online advertising campaigns has become a vivid market with more and more players in the field. Using the information about keywords similarity and relatedness, combined with prices of keywords in advertising services, such as AdWords, it is possible to find a combination of keywords that costs less, but drives the same or bigger amount of relevant traffic. Such services, however, do not take advantage of keywords that become occasionally relevant. For instance, let us consider the situation happened at this year's SXSW[7] conference held at Austin, Texas, USA. Many new iPad applications were showcasted at the conference and a rumor appeared, and lately became truth, that iPad 2 would start selling on the second day of the conference. This trend would be noticed if NMD was measured for the words 'ipad' and 'sxsw' A company selling iPad accessories, would in such an occasion have a clear interest to alter the keywords for their AdWords campaign for promoting its products and add the word 'sxsw', thus getting new relevant traffic. Once the NMD for the two words goes up again, the advertising campaign can again be changed to avoid driving the traffic that became less relevant.

Responding to changes in terms' relatedness over time, for advertising campaigns means not missing out relevant traffic, and as such is of high importance for this market. Web marketing tools such as KeywordDiscovery.com do offer the possibility to discover relevant keywords and include them in marketing campaigns, but do not reflect the change in this relevancy. Changes in relevancy might open completely new possibilities for advertising campaign optimization, and using our notion of NMD, these changes may even be taken into account in an automated or semi-automated way.

### Scenario 2: Facilitating Discovery of Relevant Resources in Organizations

Many organizations, especially larger ones, maintain organizational vocabularies and use them for the annotation of different kinds of documents and other digital assets. Such a vocabulary often results from a collaborative work of domain experts and a knowledge engineer. Therefore, it tends to reflect the experts' view of the subject

---

[7] http://sxsw.com/

domain, and the terms it defines reflect the jargon used by these experts. However, this jargon does not necessarily overlap with the everyday language used by the employees within the organization. As a consequence, employees would experience difficulties in formulating their requests for different kinds of organizational resources using the organization's official vocabulary. This indicates the need for harmonizing the official and the actual vocabularies within an organization. Furthermore, each organization evolves and many organizations need to go through continuous changes in order to respond to the constantly changing conditions in their environment. To properly address the evolving work practices in the organization, the organization's vocabulary has to evolve as well, and it should evolve to be comprehensible and usable by the employees (i.e., it should incorporate the terminology used by the employees). This is where the suggested dynamic MSR applied over the messages exchanged in the organization's Twitter-like communication channels (e.g., Yammer[8]) can help. In particular, the proposed MSR can be used for extracting terms related to certain tasks, projects, organizational positions, etc., in order to use them for evolving the organization's vocabulary. This would increase the usability of the vocabulary and consequently improve the search and discovery of organizational resources.

The suggested dynamic MSR can also be applied for facilitating people search within an organization by enabling the deduction of terms that best describe each employee. Previous studies exploring the practice of people tagging in organizations [23][24] have confirmed that people do perceive such a practice beneficial as it allows for, e.g., finding out who is working on a certain project/task, or identifying experts in a particular topic. However, the main obstacle for applying this practice in workplace lies in the very act of directly tagging (labeling) a person; many participants in the cited studies were reluctant to directly tag their colleagues as they were worried about potentially inadvertent effects those tags might cause. With the proposed dynamic MSR applied to the messages exchanged within the organization's micro-blogging and/or social streaming application, an organization would be able to identify the terms (tags) related to each employee. These terms would still reflect the community's perception of any particular employee, while freeing people from the unnecessary cognitive burden of inadvertently affecting their colleagues.

## 7  Example Diagrams

In order to test the use of the formula (1) on the data gathered from Twitter in several consecutive days for detection of the change in SR between two terms, we chose several examples of term pairs whose popularity we, as humans, were able to perceive from the news. The testing was done using our Tweet Dynamics application (cf. Section 5).

Since, unfortunately, catastrophic events were happening in Japan at the time of writing this paper[9], we used keywords 'japan' and 'nuclear' and calculated their NMDs for 5 days starting from March 8, 2011.

---

[8] https://www.yammer.com/

[9] On March 11, 2011, a strong earthquake struck Japan which triggered a failure of the cooling system of the reactor at Japan's Fukushima nuclear power plant, causing a huge explosion at the power plant the day after, on March 12.

Figure 3 - NMD diagram for terms 'japan' and 'nuclear' for the 5 days period

By looking at the diagram (Figure 3), one can observe that by March 11, there was a small relatedness between the terms 'japan' and 'nuclear' because the earthquake happened suddenly; thus the value of NMD (shown on Y axis) is higher. On the day of the earthquake (March 11[th]), one can see that the NMD significantly decreased, i.e., SR of the terms increased, as many people tweeted about the danger of explosion at the nuclear power plant. That trend continued in the following days.
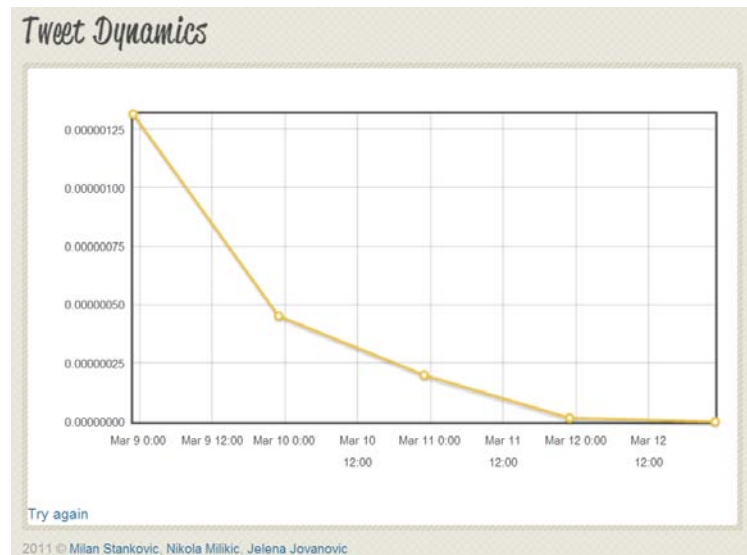


Figure 4 - NMD diagram for terms 'ipad' and 'sxsw' for the 5 days period

Our second example is about terms 'ipad' and 'sxsw' (already mentioned in Section 6). iPad started selling unexpectedly during the SXSW on the 12th of March. From the diagram, it is obvious that there was a rumor about it some days ago, as the NMD decreased exponentially towards the first day of sales, to reach its lowest value on the 12th of March. It is easy to think of potential benefits that the owners of iPad-related content might derive from this newly related term, by including it in their advertising campaigns, and using it for positioning their content. The relevant NMD diagram is shown on Figure 4.

As already mentioned, there is a big limitation of using Twitter Search API, because it limits the number of search results to a maximum of 1500. If we had access to the whole corpus of messages posted in this period, we would have been able to measure the change in relatedness more precisely. However, in the case of terms that are usually rather non-related, the importance of the change is still noticeable even with such limitations imposed.

## 8  Conclusions and Future Work

This paper presents our initial work on using data streams from Twitter or Twitter-like services for the detection of changes in semantic relatedness of terms. In particular, being inspired by the work of Cilibrasi & Vitanyi [16] on using Google search results for computing semantic relatedness of terms, we have introduced Normalized Micropost Distance (NMD). It makes use of micropost streams of Twitter-like services to compute semantic relatedness of two terms for a given time period. We have also suggested how our approach can be leveraged in two real-life scenarios that differ both in the application domain (online advertising and organizational knowledge management) and the data source to be used for the computation of the NMD measure (Twitter and organization's internal micro-blogging service).

An important challenge to attack in our future work is the detection of good candidate term pairs, i.e., pairs where a change is likely to happen. Our NMD measure allows one to measure the change in semantic relatedness, and follow it over time, but does not directly help in identifying which term pairs are likely to be the subject of change without calculating the NMD values for all possible term pairs. Having such a possibility is important in light of the need for computational efficiency and of the limits imposed by Twitter and other major players on Real-time Web. The detection of candidates for NMD calculation is dependent of the actual usage scenario, as each real-life scenario is related to a specific subject domain characterized by its specific language and important topics. Accordingly, for each scenario, there would be a list of terms to watch. With such a list available, it would be enough to identify the candidate terms that, when coupled with the watched terms could form pairs for which the calculation of NMD might lead to the detection of significant relatedness. We believe that looking at trending topics on Twitter, as well as in recent news articles, might help in finding good candidate terms for a Web marketing scenario (as presented in Section 6). Our intention is thus to explore this research question and deliver a system that could take a number of terms to watch, and provide a list of terms that have recently become more related to one or more of the watched terms.

Another equally important direction of our future work is a comprehensive evaluation of the proposed dynamic measure of semantic relatedness of terms. For that purpose we intend to use Twitter's Streaming API[10], and in particular its "Gardenhose" access level which offers the proportion of the Twitter's public data stream (currently, around 10%) that could form a statistically significant sample. This approach would help us overcome the mentioned limitations of using Twitter Search API. Besides that, since Google recently started including real–time updates coming from Twitter, it could serve us as an important source of data. But since, at the time of writing this paper, these data were not accessible through Google Search API, we need to wait for this feature to become programmatically available. Using this data stream, we intend to do an evaluation study that would consist of a comparative analysis of our approach and the approach we found as the most related to our work, namely the approach reported in Song et al. [21].

# References

[1] Wagner, C. (2010). Exploring the Wisdom of the Tweets: Towards Knowledge Acquisition from Social Awareness Streams. PhD Symposium at 7th Extended Semantic Web Conference (ESWC2010) Heraklion, Crete, Grece: Springer. Retrieved March 10, 2011, from http://www.springerlink.com/index/R4463T1333777N11.pdf.

[2] Sheth, A., Thomas, C., & Mehra, P. (2010). Continuous Semantics to Analyze Real-Time Data. IEEE Internet Computing 14, 6 (November 2010), 84-89.

[3] Stankovic, M., Rowe, M., & Laublet, P. (2010). Mapping Tweets to Conference Talks: A Goldmine for Semantics. *in Proceedings of the 3rd Social Data on the Web Conferece, SDOW2010, collcoated with International Semantic Web Conference ISWC2010.* Shanghai, China.

[4] Sigurbjörnsson, B., & Zwol, R. van. (2008). Flickr tag recommendation based on collective knowledge. Proceeding of the 17th international conference on World Wide Web - WWW '08, 327. New York, New York, USA: ACM Press. doi: 10.1145/1367497.1367542.

[5] Mei, Q., Zhou, D., & Church, K. (2008). Query suggestion using hitting time. *Proceeding of the 17th ACM conference on Information and knowledge mining - CIKM '08*, 469. New York, New York, USA: ACM Press. doi: 10.1145/1458082.1458145.

[6] Safar, B., & Kefi, H. (2004). OntoRefiner, a user query refinement interface usable for Semantic Web Portals. *Proceedings of Application of Semantic Web technologies to Web Communities, Workshop ECAI'04* (pp. 65-79). Retrieved January 25, 2011, from http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:OntoRefiner,+a+user+query+refinement+interface+usable+for+Semantic+Web+Portals#0.

[7] Macdonald, C., & Ounis, I. (2007). Expertise drift and query expansion in expert search. *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management - CIKM '07*, 341. New York, New York, USA: ACM Press. doi: 10.1145/1321440.1321490.

[8] Cross, V., "Semantic Relatedness Measures in Ontologies Using Information Content and Fuzzy Set Theory," In Proc. of the 14th IEEE Int'l Conf. on Fuzzy Systems, (2005), pp. 114–119.

[9] Gasevic, D., Zouaq, A., Torniai, C., Jovanovic, J., Hatala, M., "An Approach to Folksonomy-based Ontology Maintenance for Learning Environments," IEEE Transactions on Learning Technologies, 2011 (in press)

[10] Waltinger, U., Cramer, I., & Wandmacher, T. (2009). From Social Networks To Distributional Properties: A Comparative Study On Computing Semantic Relatedness. Cognitive Science.

[11] Burton-Jones, A., Storey, V., Sugumaran, V., & Purao, S. (2003). A heuristic-based methodology for semantic augmentation of user queries on the web. *Conceptual Modeling-ER 2003*, 476–489.

---

[10] http://dev.twitter.com/pages/streaming_api

Springer. Retrieved January 19, 2011, from http://www.springerlink.com/index/TP1URDMGDM3F0WP3.pdf.

[12] Ziegler, C.-N., Simon, K., & Lausen, G. (2006). Automatic Computation of Semantic Proximity Using Taxonomic Knowledge Categories and Subject Descriptors. *CIKM '06 Proceedings of the 15th ACM international conference on Information and knowledge management* (pp. 465-474). Arlington, Virginia, USA: ACM New York, NY, USA. Maguitman, A. G., Menczer, F., Roinestad, H., & Vespignani, A. (2005). Algorithmic detection of semantic similarity. *Proceedings of the 14th international conference on World Wide Web* (p. 107–116). ACM. Retrieved January 25, 2011

[13] Resnik, P. (1995). Using Information Content to Evaluate Semantic Similarity in a Taxonomy. Arxiv preprint cmp-lg/9511007, 1. Retrieved January 24, 2011, from http://arxiv.org/abs/cmp-lg/9511007

[14] Matos, S., Arrais, J. P., Maia-Rodrigues, J., & Oliveira, J. L. (2010). Concept-based query expansion for retrieving gene related publications from MEDLINE. *BMC bioinformatics*, *11*, 212. doi: 10.1186/1471-2105-11-212.

[15] Salton, G. and McGill, M. Introduction to Modern Information Retrieval. New York: McGraw-Hill, 1983

[16] Cilibrasi, R. L., & Vitanyi, P. M. B. (2007). The Google Similarity Distance. *IEEE Transactions on Knowledge and Data Engineering*, *19*(3), 370-383. doi: 10.1109/TKDE.2007.48.

[17] Gracia, J., & Mena, E. (2008). Web-Based Measure of Semantic Relatedness. In Proceedings of the 9th international conference on Web Information Systems Engineering (WISE '08), pp.136-150.

[18] Strube, M., & Ponzetto, S. P. (2006). WikiRelate! Computing semantic relatedness using Wikipedia. *Proceedings of the National Conference on Artificial Intelligence* (Vol. 21, p. 1419). Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999. Retrieved February 22, 2011,

[19] Mendes, P. N., Passant, A., & Kapanipathi, P. (n.d.). Twarql: Tapping into the Wisdom of the Crowd. *Proceedings of the 6th International Conference on Semantic Systems* (p. 1–3). Graz, Austria: ACM. Retrieved March 14, 2011, from http://portal.acm.org/citation.cfm?id=1839762.

[20] Sakaki, T., Okazaki, M., & Matsuo, Y. (2010). Earthquake shakes Twitter users: real-time event detection by social sensors. *Proceedings of the 19th international conference on World wide web* (p. 851–860). ACM. Retrieved March 6, 2011

[21] Song, S., Li, Q., and Zheng, N. (2010). A spatio-temporal framework for related topic search in micro-blogging. In Proceedings of the 6th international conference on Active media technology (AMT'10), Aijun An, Pawan Lingras, Sheila Petty, and Runhe Huang (Eds.). Springer-Verlag, Berlin, Heidelberg, 63-73.

[22] Nagarajan, M., Gomadam, K., Sheth, A., Ranabahu, A., Mutharaju, R., Jadhav, A.: Spatio-temporal-thematic analysis of citizen sensor data: Challenges and experiences. Web Information Systems Engineering-WISE 2009 pp. 539-553 (2009)

[23] Farrell, S., Lau, T., Wilcox, E., and Muller, M. "Socially Augmenting employee profiles with people-tagging," Proceedings of the 20th annual ACM symposium on User interface software and technology, Newport, Rhode Island, USA, 2007, pp. 91-100

[24] Braun, S., Kunzmann, C., & Schmidt, A. (2010). People Tagging & Ontology Maturing: Towards Collaborative Competence Management. In: David Randall and Pascal Salembier (eds.): From CSCW to Web2.0: European Developments in Collaborative Design, Selected Papers from COOP08, Springer, Berlin/Heidelberg.