



Evaluating authoritative sources using social networks: an insight from Wikipedia

Nikolaos Th. Korfiatis

*Department of Informatics, Copenhagen Business School (CBS),
Copenhagen, Denmark, and*

Marios Poulos and George Bokos

Department of Archives and Library Sciences, Ionian University, Corfu, Greece

Abstract

Purpose – The purpose of this paper is to present an approach to evaluating contributions in collaborative authoring environments, and in particular, Wikis using social network measures.

Design/methodology/approach – A social network model for Wikipedia has been constructed, and metrics of importance such as centrality have been defined. Data has been gathered from articles belonging to the same topic using a web crawler, in order to evaluate the outcome of the social network measures in the articles.

Findings – Finds that the question of the reliability regarding Wikipedia content is a challenging one and as Wikipedia grows, the problem becomes more demanding, especially for topics with controversial views such as politics or history.

Practical implications – It is believed that the approach presented here could be used to improve the authoritativeness of content found in Wikipedia and similar sources.

Originality/value – This work tries to develop a network approach to the evaluation of Wiki contributions, and approaches the problem of quality Wikipedia content from a social network point of view.

Keywords Social networks, Encyclopaedias

Paper type Research paper

1. Introduction and background

Since, the invention of writing as a method of encoding human knowledge, the preservation and dissemination of information and knowledge has become a matter of great importance to humanity. People, as intelligent entities, produce and consume information that is preserved in, and accessed from, various sources such as books, articles and encyclopaedias. Since, a high level of complexity characterizes the organization of information, reference works that assist the retrieval of relevant information resources are crucial for dissemination and further development of human knowledge in a particular subject. Encyclopaedias and dictionaries represent the major instances of such reference works, since their principal scope is to assist, through associative trailing, the retrieval of the relevant resources through a particular domain (collection of relevant lemmas).

Visions of the world wide web such as the Memex, envisioned by Bush (1945), and the original intuition behind the design of the world wide web by Berners-Lee *et al.* (1994), tend to represent the world wide web as a huge encyclopaedia, where lemmas



are associated by using a hypermedia model. Nonetheless, an encyclopaedia and any other kind of reference work is subject to evaluation as to the level of its quality characterizing. Since, the production of such works is subject to a very small number of individuals, the development process is characterized by high complexity. Efforts in the world wide web such as Wikipedia, try to distribute that kind of complexity to several contributing authorities by allowing the synchronous editing and publication of lemmas through its publication model. However, since its inception Wikipedia has been subject to criticism (Fasoldt, 2004; Orlowski, 2005; Lipczynska, 2005), due to what level the information contained can be trusted and referenced in research works. In that case, models of credibility, which are used extensively on search engine research and information retrieval, can be used to evaluate the trustworthiness of the topic covered by Wikipedia.

On the web, several successful approaches to credibility such as PageRank (Brin and Page, 1998; Page *et al.*, 1999) use methods derived from graph theory to model credibility, which utilize the connections of the resource for evaluation. Several graph theoretic models of credibility and text retrieval (Faloutsos, 1985) rely on the consideration of the in-degree of the node (the sum of the incoming arcs of a node in a directed graph) to extract importance and trustworthiness. This is also implied by the publication workflow and the resulting context on which those models are based. For instance, on the world wide web and similar hypermedia systems, such models of credibility evaluate a web page using the in-degree extracted by the hypertextual context. However, there are publication models supporting social activities (e.g. collaborative authoring) that derive much of their credibility from their productions (e.g. authorship), where the hyperlink context does not depict that kind of activity. In that case, the in-degree cannot provide input to evaluate the importance of that entity; therefore, a holistic approach is needed. This alternative evaluation has to consider the outputs of the entity (productions), as happens with several informal social communication models (Festinger, 1950). In a graph theoretic interpretation this can be modelled as the outer-degree of the node, which conceptually represents the entity evaluated.

2. The Wiki publication model

The web has given rise to new forms of collaboration and interaction facilitating the manipulation of shared artefacts and information spaces (Cadiz *et al.*, 2000). In the current state, the web ecosystem consists of resources (web pages/ files) linked through hypertext connectors, thus forming a system of links denoting references to those resources as well as providing views to requesting authorities.

However, one of the initial design goals of the web was not only to facilitate views of the resources requested but also allow editing and annotation of these resources in a simple way (Berners-Lee and Fischetti, 1999). The concept of Wiki (Leuf and Cunningham, 2001) has given a response to this challenge. WikiWiki (Hawaiian for “quick”) applications facilitate a way of collaborative editing, supported by a revision mechanism that allows the monitoring of changes and contributions to the sections edited. The use of WikiWiki applications is common in cases such as formation of collaborative document editing (e.g. in communities of practice), or formation of shared knowledge repositories such as dictionaries, etc. One of the best-known implementations of the usefulness of the Wiki system to support collaborative document editing is the Wiki-based encyclopaedia Wikipedia and its related projects (www.wikimedia.org).

Traditional encyclopaedias such as Britannica are often characterized by a high level of credibility by domain experts, taking into account the background process that has resulted – domain authorities contribute to the final outcome. Alternatively, since it uses the WikiWiki system, Wikipedia allows the editing and creation of encyclopaedic articles by anyone who wishes to contribute. Its primary aim is to provide free editing access and gather knowledge representing the consensus of the term presented, and thus not to evaluate the contributing authorities. However, as the content increases along with the contributing sources (Figure 1), a critical issue has emerged regarding the credibility of Wikipedia as an authoritative reference source (Andrew *et al.*, 2005; Lih, 2004). The issue is extended not only to the outcome (article) but also to the process of shaping the article, in which a contributor would allow another authority to submit, change or delete a contribution accepted or not accepted by him/her. Wikipedia has internal mechanisms of managing those cases such as a permission ranking system, where a contributor is accredited by the level of participation in the shaping of the article, as well as a discussion tab on most of the articles or notifications and warnings regarding the content. Nevertheless, the research question looks at how to provide a clue to the credibility for an article based on the contributing authorities, and their acceptance by the community of their fellow contributors.

Values on *X*-axis represent the articles on logarithmic scale. Values on *Y*-axis represent the number of contributors. (Statistics obtained from <http://en.wikipedia.org/wiki/stats/EN/TablesWikipediaEN.htm>).

In this paper we present an initial attempt to model the problem towards providing an authoritative ranking mechanism based on social interaction data collected through the Wiki. Social interaction is approached from social communication facilitated by the Wikipedia platform (e.g. edits on edits) (Cobley, 1996). We then model the credibility of each contributor using the metric of centrality, thus producing an overall centrality

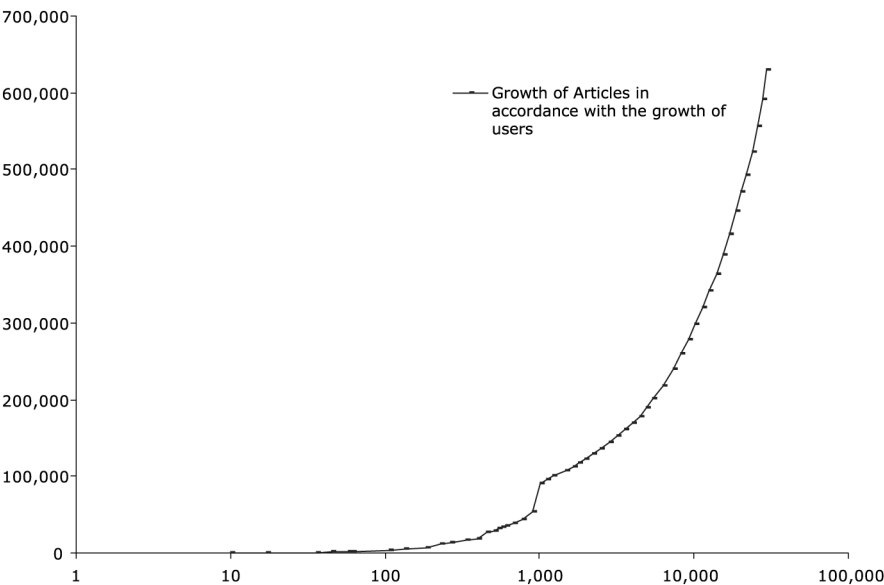


Figure 1.
Growth of articles in
accordance with the
growth of users in the
English Wikipedia

measure for the article depicting the social activity/process that has taken place through the shaping of the article. We argue that this factor can be used as a metric of credibility, representing the article and the contributing authorities.

3. A network approach on the Wiki publication model

Social networks and social network analysis in particular (Wasserman *et al.*, 1994; Scott, 2000), is a research paradigm that tries to unravel patterns of social relationships across various individuals in a social context. Following the patterns and measuring structural and compositional values in the networks, we ought to understand the basic structure and properties of the network and explain its behaviour; thus, uncover those actions that characterize most of the activity described by the network. Social network analysis focuses on a more rationalistic approach to research on organizations and social groups (Borgatti and Foster, 2003), since it aims to expand interdependent relations between the members of the group. WikiWiki applications facilitate a case where social relationships are established over a domain of social actions such as acceptance, objection or rejection of a contribution. As in the case of Wikipedia, the Wiki facilitates a collaborative document editing effort, which relies on the contribution of multiple authors in a concurrent system. This enables the combination of contributions in an effective and democratic way, allowing all the ground knowledge about the article/lemma to be present in the most recent revision of the article. By democratic, we also refer to the ability of anyone who uses the Wiki to contribute or to make modifications to content contributed by someone else. In that sense, as the Wiki-fication continues, the final document (or the most recent revision) is the outcome of a community process involving certain social interactions embedded in the content modification.

From a social research point of view, what makes such a case interesting is the negotiation process that takes place when writing and structuring the article. For example, a user makes a contribution that is erased, and the user tries to establish its contribution back (to make it visible and accepted by the others). In both cases (article and negotiation), there are interaction ties characterizing the final outcome and the dynamics of the process. In this paper, we will focus on the interaction ties between multiple contributors working on the same article or domain of articles in the Wikipedia namespace. However, to do this kind of study, we first need to define the structural and compositional variables which characterize such a network.

In the Wiki publication model, we can see the necessary structural and compositional variables important in the construction of a social network of contributors to an article or topic in the Wikipedia. Structural and compositional properties of the publication model can be found in the following use-cases:

- When a contributor edits content submitted by someone else it establishes a tie with him/her. This is depicted by an acceptance factor, which represents the percentage of the previous contributor's content that is visible after.
- Every contributor who has a single contribution or more to the article establishes a relational tie with the other content contributors. Evidence of participation in common projects strengthens this tie.

We can also link actors through two different layers of networks (Figure 2):

- (1) *The articles network.* Every article in the Wikipedia contains references to other articles as well as external references. A set of links used for classification purposes is also available in most of the active articles of the encyclopaedia. Every article represents a vertex in the article network and the internal connections between the article edges of the network.
- (2) *The contributors network.* Wikipedia is a collaborative writing effort, which means that an article has multiple contributors. We assume that a contributor establishes a relationship with another contributor if they work on the same article. In the resultant signed network, a vertex represents each contributor, and their social ties (positive or negative) are represented by an edge denoting the sequence of their social interaction.

The resultant graph is a two-mode network with two sets of entities: articles and contributors. Contributors can be either connected (belong to the same article) or interconnected (common contributions on two or more articles in the same domain). In an article domain of high credibility it is expected that more interrelations will be found, since the contributors may contribute content to more than one article, thus depicting their common interest. Therefore, the more affiliated a contributor becomes with a domain, the more interested he/she is in the article; thus representing knowledge of the domain. Let us consider, for instance, a contributor who has made a lot of contributions to the domain regarding the history of Spanish colonies in Latin America. The author has also completed some contributions to an article on anatomy. However, the author is more affiliated to the articles regarding the history of the Spanish colonies than to medicine. Therefore, his contribution to medicine may be considered less authoritative than those contributions made in the other domain, as his knowledge of anatomy is not as extensive as his knowledge of Spanish colonies.

In social network analysis, there are a variety of measures that can assess this kind of social activity in a sociometric study. As we have already defined our graph, we can use some common social network metrics to extract this kind of information from Wikipedia data.

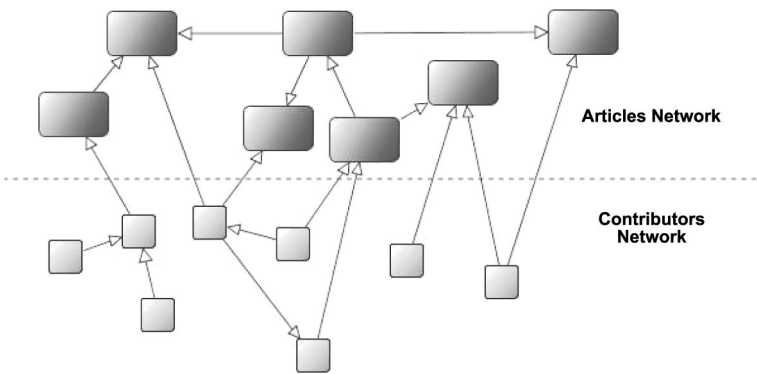


Figure 2.
Network layers in the
Wiki publication model.
Contributors are linked
together by working on
common projects (articles)
in the Wikipedia
namespace

4. Network measures in the Wikipedia contributions

As previously mentioned, contributors contribute to one or more articles belonging to the same or different domains. Based on this, we can evaluate the activity of the contributors and thus extract metrics of their presence in the domain of the articles. The development of those metrics is based on the following assumptions:

- The more decentralized the editing of an article, the better the article represents a consensus.
- The contributors whose content has been most accepted (seen from the result of the *diff operation* in the Wiki) are attributed a level of authority regarding the article.
- This level of authority remains only in the domain of the article. However, domains that belong to the same topic retain the level of authority for a contributor.

The graph we model is a signed directed network, with arcs as a factor depicting the level of acceptance of the content submitted by contributor *A* and accepted by contributor *B*. In order to model the authoritativeness of contributors, we selected the centrality index (Freeman, 1979; Sabidussi, 1966) of the resultant graph and, in particular, a measure of centrality dealing with the degree (total of incoming, outgoing edges) of the vertex/contributor in the examined article. The concept of centrality is well accepted in social network analysis, as there are numerous studies showing the usefulness of such a metric for measuring activity in social networks (Everett and Borgatti, 1999; Freeman, 1979). In sociometric studies, the usage of centrality is targeted to unfold the person/individual who is the most prominent in a network, thus ranking the actors according to their positions in the network; and is interpreted as the prominence of actors embedded in a social structure. In our study, we use the degree centrality index, which is the simplest definition of centrality and is based on the incoming and outgoing adjacent connections to other contributors in an article graph. To measure the centrality at an individual level, we define the contributor degree centrality; and to an article level, the article degree centralization, which represents the collective of contributor degree centrality.

4.1 Contributor degree centrality

In classical social network models, the inner degree (the amount of edges coming into a node) represents the choices the actor has over a set of other actors. However, in our Wiki network model, the amount of incoming edges represents edits to the text; therefore, the metric of inner degree is the opposite, meaning that the person with the biggest inner degree has the biggest amount of objection/rejection in the contributor community, and thus receives a kind of negative evaluation from his/her fellow contributors. On the other hand, the outer-degree of the vertex represents edits/participation in several parts of the article, and thus gives a clue to the activity of the person in relation to the article and the domain. Mathematically, we can represent such formalism as follows: considering a graph representing the network of contributors for an article contributed in the Wiki, the contributor degree centrality – a contextualized expression of actor degree centrality – is a degree index of the adjacent connections between the contributor and others who edit the article. From graph theory, the outer degree of a vertex is the cumulative value of its adjacent connections:

$$C_D(n_i) = d(n_i) \sum_j x_{ij}$$

The adjacent x_{ij} represents the relational tie between the contributors and their contribution over the domain of the article. This also is characterized by the visibility of the contribution in the final article and can be either 1 or 0. To provide the centrality, we divide the degree with the highest obtained degree from the graph, which in graph theory is proved to be the number of remaining vertices (g) minus the self ($g-1$). Therefore, the contributor degree centrality can be calculated as:

$$C'_D(n_i) = \frac{d(n_i)}{g-1}$$

4.2 Article degree centralization

We define an article's degree centrality C_{DM} as the variability of the individual contributor centrality indices. The $C_D(n^*)$ represents the largest observed contributor degree centrality:

$$C_{DM} = \frac{\sum_i^g [C_D(n^*) - C_D(n_i)]}{(g-1)(g-2)}$$

Again we divide the variability with the highest variability observed in the graph. Having defined the metrics, we apply them and explain their qualitative values in a case study of the English language Wikipedia.

5. An insight from Wikipedia

As previously mentioned, Wikipedia follows a hypermedia model to categorize the articles (lemmas) in an associative taxonomy. In that particular taxonomic classification, we define the following structures:

- *Domain*. A collection of articles that tackle a common subject (e.g. philosophy).
- *Category*. A collection of domains that have a common categorical and etymological root. For example, the domains philosophy and economics have a connection in the category of social sciences.

In order to provide a qualitative analysis of the metrics deployed in the evaluation of the articles, we picked ten articles with a similar number of contributors to the domain philosophy from the English language Wikipedia. Table I shows the list of articles used in the case study, as well as the values of their article degree centralization. The data for each article was collected using the python Wikipedia robot framework (www.pywikipedia.org). The resulting networks contained an average of 259 contributors per article, and the average article inter-relations per contributor were approximately two. We used the diff function of the Wiki to assess the tie between the pair of contributors as modelled in Section 2. The data was then analysed using a python script, in order to calculate the individual contributor degree centrality along with the article degree centralization.

Table I.
The Wikipedia articles
from which empirical
data was gathered

Article name	Number of contributors	Article degree centrality (max 1)
Adam Smith	276	0.039114
Aristotle	274	0.0232
Immanuel Kant	231	0.20484
Johann Wolfgang von Goethe	242	0.016682
John Locke	292	0.008581
Karl Marx	232	0.006601
Ludwig Wittgenstein	220	0.006328
Philosophy	280	0.00254
Plato	284	0.001207
Socrates	289	0.000405

As can be observed from the table, the article degree centralization is relatively low because of the small collections of articles used in the case study and the inter-connections of the actors in the domain. However, it is enough to let us discuss some qualitative interpretations:

- The dispersion of the actor indices denotes how dependent this article is on individual contributors. For instance, if an article has a very low degree of centralization, then it means that the social process to shape it was highly distributed, thus resulting in an article that has been submitted by multiple authorities. In our case, the articles represent a low degree of centralization, which means that contributions have been made by individuals who have interests in other domains as well.
- The range of the group degree centralization reflects the heterogeneity of the authoring sources of the article. In our case, the article “Immanuel Kant” has a significantly higher degree of centralization (Figures 3 and 4 and Table II), which means that it has been contributed to by authorities concentrated in the domain of the article, and thus who have contributed to other articles.

Contributors with higher inter-relation over the same domain represent higher authorities, based on the assumption that their interest spans the domain to which the article belongs and, therefore, have conducted background research regarding the material they have contributed.

On the other hand, contributors with lesser authority tend to have their content erased/objected by contributors with higher authority. As can be observed from Figure 3, there exist a number of contributors subject to objections regarding their submissions, and therefore, they are situated on the periphery, whereas contributors with accepted contributions (authorities) tend to be in the centre.

6. Discussion and further research

The question of the reliability regarding Wikipedia content is a challenging one. As long as the size of Wikipedia grows, the problem becomes more demanding, especially for topics with controversial views such as politics or history. Our study represents an early attempt at getting to the problem and thus working towards a more sophisticated solution to address it. However, there are a number of open issues that can extend the merit of this report.

Figure 3.
Visualization of the social network of the contributors for the article “Immanuel Kant”. Nodes in the core denote high degree centrality

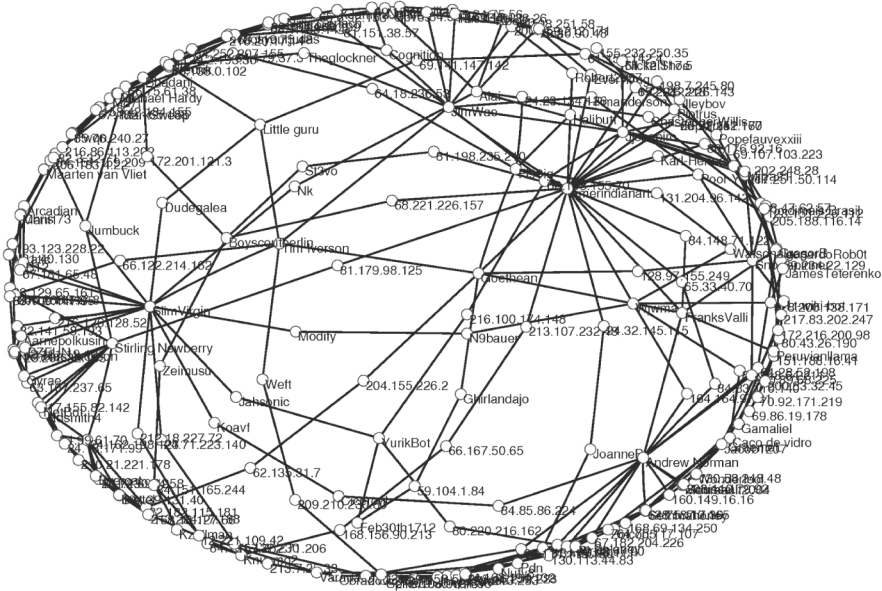
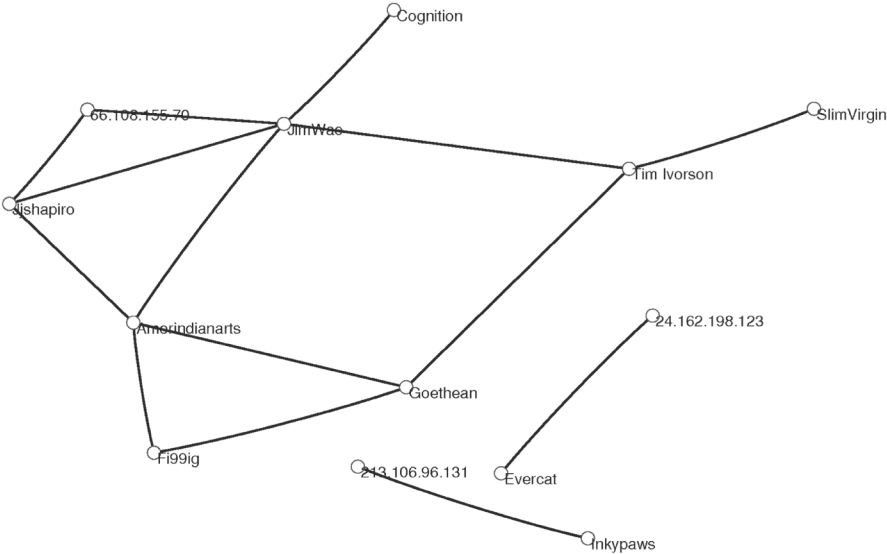


Figure 4.
A decomposition of the network to the contributors with the highest degree centrality for the article “Immanuel Kant”



The in-degree can be calculated using a more sophisticated factor, representing how much of the text contributed by one actor has been edited by another. In our case, we represent the editing or the objection by using a scale from 0 to 1, thus aggregating the factors using simple sums. A fuzzy operator could provide a solution for aggregating the results obtained by undertaking a fuzzy diff between the current version of the

Cluster (outerdegree)	Freq	Freq per cent	CumFreq	CumFreq per cent	Representative
1	1	0.4329	1	0.4329	65.6.92.153
2	199	86.1472	200	86.5801	82.3.32.71
4	14	6.0606	214	92.6407	80.202.248.28
6	6	2.5974	220	95.2381	Snowspinner
8	3	1.2987	223	96.5368	Tim Ivorson
10	1	0.4329	224	96.9697	StirlingNewberry
12	2	0.8658	226	97.8355	24.162.198.123
16	2	0.8658	228	98.7013	JimWae
18	1	0.4329	229	99.1342	Jjshapiro
20	1	0.4329	230	99.5671	SlimVirgin
31	1	0.4329 s	231	100	Amerindianart

Table II.
Contributor degree
centrality for the article
“Immanuel Kant”

article and the version submitted. In that case, the social tie also needs to be expressed in terms of fuzziness, along with the relevant cases. Expressions of credibility using imprecise criteria (Sicilia and García, 2004, 2005) can also contribute to further advancement in that direction.

The organization of topics and the definition of inter-connections is also a matter for research, since there are related domains with contributing authorities. For instance, in the category of the social sciences, a contributor who edits the article of Adam Smith and has an acceptance factor can be retained on both the economics and philosophy domains, as an article about Adam Smith is represented in both. In that case, network modelling using two layer networks (document reference, authority reference) can enhance the trust of the contributions (Hess *et al.*, 2006).

Furthermore, the measures developed and presented in this report do not actually measure the subjective quality of an article, since such a task is a cognitive process characterized by a high level of complexity. Those measures can contribute to providing an indicator of consensus related to an article, and thus assert that it does not provide controversial views or expressions of a small group of persons (especially in articles with political content). Thus, a level of neutrality expressed in the writing of the article is asserted.

Finally, specific attention should be given to the diffusion of different affiliations related to one actor. For example, a contributor may have many affiliations to unrelated subjects. This may imply that the contributor has knowledge of both fields, but in special topic cases (e.g. cardiology), the contribution in subjects such as Renaissance can be attributed as a non-expert one. Therefore, a classification of the competencies of each contributor may need to be promoted to strengthen their credibility and association with the subject or the article contributed.

References

- Andrew, L., Jakob, V., Cathy, M., Samuel, K. and Reinhold, H. (Eds) (2005), *Proceedings of Wikimania 2005 – The First International Wikimedia Conference*.
- Berners-Lee, T. and Fischetti, M. (1999), *Weaving the Web: The Original Design and Ultimate Destiny of the World Wide Web by its Inventor*, Harper, San Francisco, CA.
- Berners-Lee, T., Masinter, L. and McCahill, M. (1994), Uniform Resource Locators (URL), RFC 1738.

- Borgatti, S.P. and Foster, P.C. (2003), "The network paradigm in organizational research: a review and typology", *Journal of Management*, Vol. 29 No. 6, pp. 991-1013.
- Brin, S. and Page, L. (1998), "The anatomy of a large-scale hypertextual web search engine", *Proceedings of the Seventh International Conference on World Wide Web*, available at: www-db.stanford.edu/%7Ebackrub/google.html, pp. 107-17.
- Bush, V. (1945), "As we may think", *The Atlantic Monthly*, July.
- Cadiz, J.J., Gupta, A. and Grudin, J. (2000), "Using web annotations for asynchronous collaboration around documents CSCW '00", *Proceedings of the 2000 ACM Conference on Computer Supported Cooperative Work*, pp. 309-18.
- Cobley, P. (1996), *The Communication Theory Reader*, Routledge, London.
- Everett, M.G. and Borgatti, S.P. (1999), "The centrality of groups and classes", *Journal of Mathematical Sociology*, Vol. 23 No. 3, pp. 181-201.
- Fasoldt, A. (2004), "Librarian: don't use Wikipedia as a source", *Syracuse Post Standard*, 25 August.
- Faloutsos, C. (1985), "Access methods for text", *ACM Computing Surveys*, Vol. 17 No. 1, pp. 49-74.
- Freeman, L.C. (1979), "Centrality in social networks: conceptual clarification", *Social Networks*, Vol. 1 No. 3, pp. 2152-39.
- Festinger, L. (1950), "Informal social communication", *Psychological Review*, Vol. 57 No. 5, pp. 271-82.
- Hess, C., Stein, K. and Schlieder, C. (2006), "Trust-enhanced visibility for personalized document recommendations", *Proceedings of the 21st ACM Symposium on Applied Computing, Dijon, France*.
- Leuf, B. and Cunningham, W. (2001), *The Wiki Way: Collaboration and Sharing on the Internet*, Addison-Wesley, Reading, MA.
- Lih, A. (2004), "Wikipedia as participatory journalism: reliable sources? Metrics for evaluating collaborative media as a news resource", *Proceedings of the International Symposium on Online Journalism*.
- Lipczynska, S. (2005), "Power to the people: the case for Wikipedia", *Reference Reviews*, Vol. 19 No. 2.
- Orlowski, A. (2005), "Wikipedia science 31% more cronky than Britannica's", *The Register*.
- Page, L., Brin, S., Motwani, R. and Winograd, T. (1999), "The PageRank citation ranking: bringing order to the web", *Technical Report*, Stanford Digital Libraries Project.
- Sabidussi, G. (1966), "The centrality index of a graph", *Psychometrika*, Vol. 31, pp. 581-603.
- Scott, J. (2000), *Social Network Analysis: A Handbook*, 2nd ed., Sage, London.
- Sicilia, M.A. and Garcia, E. (2004), "Fuzzy group models for adaptation in cooperative information retrieval contexts", *Lecture Notes in Computer Science 2932*, Springer, New York, NY, pp. 324-34.
- Sicilia, M.A. and Garcia, E. (2005), "Filtering information with imprecise social criteria: a FOAF-based backlink model", *Proceedings of the Fourth Conference of the European Society For Fuzzy Logic and Technology (EUSLAT)*.
- Wasserman, S., Faust, K. and Iacobucci, D. (1994), *Social Network Analysis: Methods and Applications (Structural Analysis in the Social Sciences)*, Cambridge University Press, Cambridge.