# Chapter 14 Identifying Polarized Wikipedia Articles

**Nikos Kirtsis** Patras University, Greece Jeries Besharat Patras University, Greece

**Paraskevi Tzekou** *Patras University, Greece*  **Sofia Stamou** Patras University, Greece & Ionian University, Greece

## ABSTRACT

Wikipedia is one of the most successful worldwide collaborative efforts to put together user-generated content in a meaningfully organized and intuitive manner. Currently, Wikipedia hosts millions of articles on a variety of topics, supplied by thousands of contributors. A critical factor in Wikipedia's success is its open nature, which enables everyone to edit, revise, and/or question (via talk pages) the article contents. Considering the phenomenal growth of Wikipedia and the lack of a peer review process for its contents, it becomes evident that both editors and administrators have difficulty in validating its quality on a systematic and coordinated basis. This difficulty has motivated several research works on how to assess the quality of Wikipedia articles. In this chapter, the authors propose the exploitation of a novel indicator for the Wikipedia articles' quality, namely information credibility. In this respect, the authors describe a method that captures the polarized (i.e., biased) information across the article contents in an attempt to infer the amount of credible (i.e., objective) information about its topic is more credible and of better quality compared to an article that discusses the editors' (subjective) opinions on that topic.

### INTRODUCTION

Wikipedia is one of the most popular social media websites that enables users to create content in a collaborative manner. As Wikipedia increases in both size and popularity, there is an urgent need to come up with effective quality assessment methods that would guarantee the value of its contents. Such need is primarily imposed by Wikipedia's open nature, which enables everyone contribute new or modify existing content on a variety of topics, without any pre-requisite that

DOI: 10.4018/978-1-4666-2494-8.ch014

content insertions and/or modifications undergo a peer review process. Wikipedia's open nature has led to its remarkable growth, but at the same time, it has raised skepticism about the quality of its contents, considering that anyone can become a Wikipedia editor. In light of the above, numerous researchers over the last few years attempted to design methods and techniques that would capture the article features that signify quality and thus be able to quantify the overall quality of Wikipedia (Stvilia, et al., 2005b; Blumenstock, 2008a).

Most of existing Wikipedia quality assessment efforts, estimate the articles' value based on the study of their internal characteristics such as their contextual elements (Stvilia, et al., 2005a), their linkage in the Wikipedia graph (Kamps & Koolen, 2009), their length (Blumenstock, 2008b), their factual accuracy (Giles, 2005), the formality of their language (Emigh & Herring, 2005), and many more. Additionally, over the last couple of years researchers proposed methods for the automatic identification of controversial or vandalized Wikipedia articles (Vuong, et al., 2009; Potthast, et al., 2008) in an attempt to alleviate administrators from the laborious process of manually removing malicious content from the Wikipedia collection and at the same time assist readers discriminate between commonly accepted and disputed content.

In this chapter, we build upon existing Wikipedia quality assessment efforts and propose a novel method for automatically identifying articles that need undergo revisions and/or repair in order for their contents to reach good quality levels. Our method applies text mining and lexical analysis to the Wikipedia article contents in order to firstly capture highly polarized content in the articles' body and therefore deduce the credibility of the information that Wikipedia articles communicate to readers. In our work, the distinction between credible and polarized articles is defined as follows. A credible article is one that contains unbiased and objective information on the topic being discussed, whereas a polarized article is one that presents the personal viewpoints of its editors about the topic under discussion. Considering that Wikipedia is more than a Web 2.0 information source, we believe that the contents of its hosting articles should communicate reliable and solid information and not serve as a forum of misleading and disputable content. Therefore, via the exploitation of our proposed technique we aspire to assist Wikipedia administrators detect articles of subjective content and either repair or flag them as polarized.

In brief, our proposed technique operates as follows. Given a set of Wikipedia articles, we process their contents in order to discriminate between articles of objective and articles of subjective content. Note that articles of objective content are perceived as communicators of credible information, while articles of subjective content are perceived as communicators of biased information. Articles' processing is applied at two distinct yet complementary levels: lexical and semantic. Lexical processing of the article contents concerns parsing the articles to remove markup and then apply tokenization, part-ofspeech tagging, and lemmatization to the article textual elements. Having represented the article contents into canonical sequences of word tokens, we proceed with the identification of the article tokens that might express personal opinions. To that end, we explore the words' semantic orientation<sup>1</sup> as introduced in Hatzivassiloglou and McKeown (1997); a factor used to discriminate between words of positive and negative sentiments. The main idea is that not all words in a text serve as good indicators of the text's polarity. Thus, we rely on the observations of Hatzivassiloglou and McKeown (1997) and Turney and Littman (2003) that people use adjectives and adverbs to evaluate a topic or verbalize an opinion and try to infer the polarity of the Wikipedia articles as follows. Based on the article sentences that contain adjectives and/or adverbs, we apply shallow syntactic parsing and where needed anaphora resolution to their contents in order to identify whether the later refer to the topics being discussed in their

11 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/identifying-polarized-wikipedia-articles/71860

## **Related Content**

#### Efficient String Matching Algorithm for Searching Large DNA and Binary Texts

Abdulrakeeb M. Al-Ssulami, Hassan Mathkourand Mohammed Amer Arafah (2017). *International Journal on Semantic Web and Information Systems (pp. 198-220).* www.irma-international.org/article/efficient-string-matching-algorithm-for-searching-large-dna-and-binary-texts/189771

#### A Layered Model for Building Ontology Translation Systems

Oscar Corchoand Asunción Gómez-Pérez (2007). *Semantic Web-Based Information Systems: State-ofthe-Art Applications (pp. 161-189).* www.irma-international.org/chapter/layered-model-building-ontology-translation/28913

#### A Review of Fuzzy Models for the Semantic Web

Hailong Wang (2009). *The Semantic Web for Knowledge and Data Management (pp. 23-37).* www.irma-international.org/chapter/review-fuzzy-models-semantic-web/30384

#### A Software Modeling Approach to Ontology Design via Extensions to ODM and OWL

Rishi Kanth Saripalle, Steven A. Demurjian, Alberto De la Rosa Algarínand Michael Blechner (2013). International Journal on Semantic Web and Information Systems (pp. 62-97). www.irma-international.org/article/a-software-modeling-approach-to-ontology-design-via-extensions-to-odm-andowl/94599

#### Sharing Resources through Ontology Alignments in a Semantic Peer-to-Peer System

Jérôme Euzenat, Onyeari Mbanefoand Arun Sharma (2010). *Cases on Semantic Interoperability for Information Systems Integration: Practices and Applications (pp. 107-126).* www.irma-international.org/chapter/sharing-resources-through-ontology-alignments/38041