



# The added value of auxiliary data in sentiment analysis of Facebook posts



Matthijs Meire<sup>a</sup>, Michel Ballings<sup>b,\*</sup>, Dirk Van den Poel<sup>a</sup>

<sup>a</sup>Department of Marketing, Ghent University, Tweeckerkenstraat 2, Ghent, 9000 Belgium

<sup>b</sup>Department of Business Analytics and Statistics, The University of Tennessee, 249 Stokely Management Center, 916 Volunteer Blvd, Knoxville, 37996 TN, USA

## ARTICLE INFO

### Article history:

Received 10 September 2015  
 Received in revised form 17 June 2016  
 Accepted 19 June 2016  
 Available online 27 June 2016

### Keywords:

Facebook  
 Text mining  
 Sentiment analysis  
 Machine learning  
 Social media

## ABSTRACT

The purpose of this study is to (1) assess the added value of information available before (i.e., leading) and after (i.e., lagging) the focal post's creation time in sentiment analysis of Facebook posts, (2) determine which predictors are most important, and (3) investigate the relationship between top predictors and sentiment. We build a sentiment prediction model, including leading information, lagging information, and traditional post variables. We benchmark Random Forest and Support Vector Machines using five times twofold cross-validation. The results indicate that both leading and lagging information increase the model's predictive performance. The most important predictors include the number of uppercase letters, the number of likes and the number of negative comments. A higher number of uppercase letters and likes increases the likelihood of a positive post, while a higher number of comments increases the likelihood of a negative post. The main contribution of this study is that it is the first to assess the added value of leading and lagging information in the context of sentiment analysis.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

In the beginning of the century, Web 2.0 emerged as an ideological and technical foundation giving rise to the massive production of user generated-content (UGC). Blogging platforms and online retailers are the first examples of this foundation [50]. Today, UGC is still growing rapidly, sparking interest and activity in opinion mining and sentiment analysis [62, 74]. Sentiment analysis is defined as the computational process of extracting sentiment from text [61, 74]. Applications range from the prediction of election outcomes [17, 92], to relating public mood to socio-economic variables [17], to improved e-learning strategies [72].

Early examples of sentiment analysis were mainly based on review data. This type of data rarely contained much more information than the content and the time of posting of the review itself. Models using these data are based on present information, where 'present' refers to the time of posting. This changed with the advent of social networks such as Facebook and Twitter in that much more data became available. On these platforms, not only the focal post's content is available, but, taking into account the time of posting, there is also leading and lagging information. Leading information is available even before content is posted (e.g., user profiles, previous

posts) and thus contains information about the past. On the other hand, lagging information is generated a posteriori, after the content was posted (e.g., interactions such as likes or retweets) and thus contains information about the future (seen from the time of posting). Leading information can therefore be included in any sentiment model, while lagging information can be included in tools that do not require real-time sentiment analysis. To the best of our knowledge, there is no study that includes leading and lagging information into sentiment analysis models. However, we believe that we can improve sentiment prediction by including leading and lagging information for several reasons. First, social media suffer from a lot of slang [41, 72] making it harder for traditional methods to achieve satisfactory model performance on text variables alone. Second, leading variables would take into account users' average sentiment, word use, well-being, and mood and demographics, effectively acting as a user-specific informative prior of future sentiment and accounting for heterogeneity among users. Leading variables have been shown to lead to better predictions [10]. Third, extant literature has found significant relationships between post sentiment and lagging information such as likes and comments [87].

To fill this gap in literature, we assess the additional value for sentiment analysis of leading and lagging information over and above information extracted from the focal post. We do this by constructing three models. The first model is the base model that focuses on the present and contains only the focal post (including text and timing of posting). The second model contains both the focal post's content and leading information, and thus contains both present and past

\* Corresponding author.

E-mail addresses: [Matthijs.Meire@UGent.be](mailto:Matthijs.Meire@UGent.be) (M. Meire), [Michel.Ballings@utk.edu](mailto:Michel.Ballings@utk.edu) (M. Ballings), [Dirk.VandenPoel@UGent.be](mailto:Dirk.VandenPoel@UGent.be) (D. Van den Poel).



Download English Version:

<https://daneshyari.com/en/article/551974>

Download Persian Version:

<https://daneshyari.com/article/551974>

[Daneshyari.com](https://daneshyari.com)