Research article

# Opinion polarity detection in Twitter data combining shrinkage regression and topic modeling

Hyui Geon Yoon [a], Hyungjun Kim [a], Chang Ouk Kim [a,*], Min Song [b]

[a] *Department of Industrial Engineering, Yonsei University, Seoul, Republic of Korea*
[b] *Department of Library and Information Science, Yonsei University, Seoul, Republic of Korea*

A B S T R A C T

We propose a method to analyze public opinion about political issues online by automatically detecting polarity in Twitter data. Previous studies have focused on the polarity classification of individual tweets. However, to understand the direction of public opinion on a political issue, it is important to analyze the degree of polarity on the major topics at the center of the discussion in addition to the individual tweets. The first stage of the proposed method detects polarity in tweets using the Lasso and Ridge models of shrinkage regression. The models are beneficial in that the regression results provide sentiment scores for the terms that appear in tweets. The second stage identifies the major topics via a latent Dirichlet analysis (LDA) topic model and estimates the degree of polarity on the LDA topics using term sentiment scores. To the best of our knowledge, our study is the first to predict the polarities of public opinion on topics in this manner. We conducted an experiment on a mayoral election in Seoul, South Korea and compared the total detection accuracy of the regression models with five support vector machine (SVM) models with different numbers of input terms selected by a feature selection algorithm. The results indicated that the performance of the Ridge model was approximately 7% higher on average than that of the SVM models. Additionally, the degree of polarity on the LDA topics estimated using the proposed method was compared with actual public opinion responses. The results showed that the polarity detection accuracy of the Lasso model was 83%, indicating that the proposed method was valid in most cases.

© 2016 Elsevier Ltd. All rights reserved.

## 1. Introduction

Social network service (SNS) provides a new platform on which users can freely express their opinions. Users share their opinions on a variety of issues through SNS platforms, and various researchers have conducted analyses focusing on the "exchange of opinions" that occurs on SNS platforms. Twitter is a representative SNS platform. Twitter's most prominent feature is that users worldwide can communicate quickly in 140 or fewer characters (Jansen, Zhang, Sobel, & Chowdury, 2009), and features such as "following", "followers", and "retweets" enable the rapid exchange of information. Some have argued that Twitter played a role in the election of U.S. President Barack Obama (Cogburn & Espinoza-Vasquez, 2011). Furthermore, Larsson and Moe (2012) studied the effect of Twitter on Swedish election campaigns, and Graham, Broersma, Hazelhoff, and van't Haar (2013), investigated how Twitter was used by candidates in the U.K. general election.

---

* Corresponding author.
  *E-mail address:* kimco@yonsei.ac.kr (C.O. Kim).

Existing polarity (or sentiment) analysis focuses on classifying the polarity of individual texts (e.g., web reviews or tweets) by selecting important features through methods such as n-grams (Kennedy & Inkpen, 2006; Kouloumpis, Wilson, & Moore, 2011; Pak & Paroubek, 2010; Pang, Lee, & Vaithyanathan, 2002), word subsequence (Matsumoto, Takamura, & Okumura, 2005; Xia, Zong, & Li, 2011), information gain (Forman, 2003; Zhang, Ye, Zhang, & Li, 2011), and recursive feature elimination (Abbasi, Chen, & Salem, 2008; Abbasi, Chen, Thoms, & Fu, 2008). A tweet is then classified via algorithms, such as the naïve Bayes (Xia et al., 2011), maximum entropy (Pang et al., 2002; Xia et al., 2011), or support vector machine (SVM) algorithms (Abbasi, Chen, Salem, 2008; Dang, Zhang, & Chen, 2010; Moraes, Valiati, & Neto, 2013; Saleh, Martín-Valdivia, Montejo-Ráez, & Ureña-López, 2011). However, to understand the direction of public opinion on a political issue, it is important to analyze the degree of polarity on the major topics at the center of the discussion in addition to the individual tweets. For example, it is nearly impossible to identify all of the individual opinions of the public regarding a policy presented by a candidate during an election.

This study proposes a method to analyze public opinion regarding a specific political issue (or event) by utilizing polarity analysis of Twitter data at the individual tweet level as well as at the topic level. It is expected that when properly used, polarity analysis of social media can provide several benefits, including (1) reducing the cost of traditionally conducted public opinion polls, (2) augmenting such public opinion polls, and (3) producing a more accurate gauge of public opinion overall. In addition, analyzing whether a policy proposed by a candidate is viewed favorably can help candidates develop future policies.

The proposed method consists of two stages. The first stage detects polarity in tweets using shrinkage types of multivariate regression models, namely, Lasso and Ridge (Tibshirani, 1996). We define the status of an opinion as positive or negative, excluding a neutral opinion. Twitter data reflect opinions regarding candidates that were confirmed directly, and supervised learning was conducted based on these data. The dependent variable for the polarity response is set to one when support is expressed for a certain candidate and zero when an opposing opinion is introduced; furthermore, each term in tweets is set as an independent variable to examine the effect of each term on tweet polarity. The regression coefficients of the Lasso and Ridge models indicate the degree of the term's impact on the polarity response, and we consider the regression coefficients as the sentiment scores for the terms. Lasso and Ridge added a penalty section to perform a variable reduction that maximizes detection performance—the terms that are determined to be influential on the response have non-zero regression coefficients, whereas insignificant terms have coefficients near zero. Compared with SVMs (Cortes & Vapnik, 1995), which are frequently used for polarity analyses, the regression model is able to estimate the sentiment scores of the terms in the tweets.

Although the polarity of individual tweets can be detected, the results only reveal the number of opinionated tweets; they do not show what political topics people express their opinion on. Policy makers want to know the public responses to a political issue. Therefore, in addition to the polarity detection of individual tweets, it is important to capture the major topics on which public opinion is expressed and detect the polarity of such opinions. Thus, the second stage assigns a polarity code to the major topics under discussion rather than the individual tweets. In this step, a latent Dirichlet analysis (LDA) model (Blei, Ng, & Jordan, 2003) is applied to identify the major topics. LDA represents a topic as a group of terms that frequently occur in the texts under analysis. For each topic, polarity is estimated via the sentiment scores of the terms included in the group. To the best of our knowledge, our research is the first to estimate the polarities of public opinion derived from LDA topics in this manner.

An experiment using the proposed method was conducted on the 2014 mayoral election in Seoul, South Korea, which occurred on June 4, 2014. We compared the total polarity detection accuracy of the regression models with five SVM models with different numbers of input terms selected by a feature selection algorithm. The results indicated that the performance of the Ridge model was better than that of the SVM models. Additionally, the polarity of LDA topics estimated using the proposed method was compared with actual public opinion responses. The results showed that the polarity detection of the Lasso method was valid in most cases.

The remainder of this paper is organized as follows. Section 2 reviews the prior studies related to opinion polarity analysis; Section 3 details the proposed method, which includes a data-preprocessing step; Section 4 introduces the Twitter data related to a mayoral election in Korea and presents the analysis results; and Section 5 presents the conclusion of this study and discusses directions for future research.

## 2. Related works

Communicated ideas and intentions are captured by language, and we use various expressions to reveal such ideas and intentions. The basic units that constitute this expression are terms, or morphemes, and it is through a collection of these concepts that various sentences are formed and further developed to express opinions (Na, Lee, Nam, & Lee, 2009; Zhang & Liu, 2011). Even in the same context, the choice of terms expresses differences in opinions regarding a specific issue or event. Thus, it is reasonable to hypothesize that the polarity of a term can determine the polarity of the speaker's opinion regarding a specific issue. It is presumed that when negative terms (i.e., terms with negative polarizing characteristics) are used to describe an opinion regarding a specific issue, this usage can polarize the speaker's opinion of this issue in a negative direction (Breck, Choi, & Cardie, 2007).

Early research related to the field of opinion dynamics proposed models of simple opinion evolution processes, beginning with the Ising model, which is based on statistical thermodynamics and quantum mechanics spin concepts (Galam, Gefen,