# Contextual sentiment analysis for social media genres

Aminu Muhammad*, Nirmalie Wiratunga, Robert Lothian

*School of Computing Science and Digital Media, Robert Gordon University, Aberdeen, Scotland*

A B S T R A C T

The lexicon-based approaches to opinion mining involve the extraction of term polarities from sentiment lexicons and the aggregation of such scores to predict the overall sentiment of a piece of text. It is typically preferred where sentiment labelled data is difficult to obtain or algorithm robustness across different domains is essential. A major challenge for this approach is accounting for the semantic gap between prior polarities of terms captured by a lexicon and the terms' polarities in a specific context (contextual polarity). This is further exacerbated by the fact that a term's contextual polarity also depends on domains or genres in which it appears. In this paper, we introduce SMARTSA, a lexicon-based sentiment classification system for social media genres which integrates strategies to capture contextual polarity from two perspectives: the interaction of terms with their textual neighbourhood (local context) and text genre (global context). We introduce an approach to hybridise a general purpose lexicon, SentiWordNet, with genre-specific vocabulary and sentiment. Evaluation results from diverse social media show that our strategies to account for local and global contexts significantly improve sentiment classification, and are complementary in combination. Our system also performed significantly better than a state-of-the-art sentiment classification system for social media, SENTISTRENGTH.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

Sentiment analysis concerns the study of opinions expressed in text. The task of sentiment analysis comprises of the extraction of opinion polarity (positive or negative), the target or specific aspects of the target to which the opinion refers, the holder of the opinion, and the time at which the opinion was expressed [1]. Aggregation of sentiment polarity scores from a resource such as a sentiment lexicon is typically used to classify opinionated text into sentiment classes. As a result, several general purpose sentiment lexicons have been developed and made public for research, e.g., General Inquirer [2], Opinion Lexicon [3] and SentiWordNet (SWN) [4]. However, the performance of lexicon-based sentiment analysis still remains below acceptable levels. This is because the polarity with which a sentiment-bearing term appears in text (i.e. contextual polarity) can be different from its prior polarity offered by a lexicon. Two forms of semantic difference seem to contribute to this semantic gap. First, the difference in *local context*, arising from the interaction of the term with its textual neighbourhood. For example, the prior polarity of 'good' is positive, however, such polarity is changed in 'not good'. Second, the difference in *global context* arising from the difference in the typical sentiment polarity of a term

captured by a lexicon and the term's domain- or genre-specific polarity. For example, in the text 'the movie sucks', although the term 'sucks' seems highly sentiment-bearing, this may not be reflected in a general purpose sentiment lexicon. Also, as sentiment lexicons are static resources, they need to be equipped with a strategy to adapt to changing vocabulary and sentiment over time - a characteristic of social media.

In this paper, we propose an approach to account for local and global contexts in social media genres. First, we introduce strategies to account for sentiment modifiers: negations, intensifiers/diminishers, and discourse structures. Here, we leverage the fine-grained sentiment information offered by SWN. To account for discourse structures, we introduce heuristic-based discourse parsing and weighting based on the Rhetorical Structure Theory (RST) [5]. RST posits that text can be broken into non-overlapping spans in a tree-like structure with relations that may exist between any two adjacent spans. Each text span can either have the status of the central focal point of the writer's message (i.e. nucleus) or a supporting message that helps in understanding the nucleus (i.e. satellite). As our approach is heuristic-based, we avoid the need for parsers trained with text untypical of social media, yet maintain the theoretical framework of RST. Our strategies to account for local context also incorporate non-lexical modifiers commonly used to express or emphasise sentiment in social media: capitalisation, sequence of repeated character, and emoticons. Second, we introduce an approach to hybridise general purpose lexicons with

genre-specific sentiment polarities (global context) and vocabulary. The main contributions of this paper are as follows:

- We introduce a set of strategies relevant to both the social media and a high-coverage lexicon (SWN) that adjusts term prior polarity based on local context. These include strategies for negation, intensification/diminishing, discourse structure, and non-lexical modifiers.
- We introduce a strategy to adapt a lexicon to a domain by facilitating genre-specific vocabulary enhancement using distant-supervised learning.
- We provide a comparative analysis with state-of-the-art systems.

To the best of our knowledge, this is the first time SWN, together with the proposed contextual analysis are applied to sentiment classification of social media. The rest of the paper is organised as follows. Related work is presented next in Section 2, followed by our system (SmartSA) in Section 3. Evaluation results are presented and discussed in Section 4, followed by conclusions and future work in Section 5.

## 2. Related work

The task of sentiment classification involves the labelling of text with sentiment class. Several methods have been employed for the task, drawing from both supervised/unsupervised machine learning and lexicon-based unsupervised strategies. Inspired by the field of topic-based text classification, supervised methods make use of machine learning algorithms trained with sentiment-labelled data to predict sentiment class of unlabelled test documents. Although this method was shown to work well in sentiment classification, it becomes problematic when reliable and sufficient training data are difficult to obtain. This is particularly the case for the non-review-based social media where content is not associated with ratings that could be exploited as "noisy" labels. A solution to the problem of labelled data acquisition is the use of unsupervised topic modelling approaches. These typically involve the use of probabilistic topic detection methods to detect both topic and sentiment from a collection of unlabelled documents.

Machine learning sentiment classifiers tend to be highly domain/genre specific, performing well on the domain/genre of training but poorly on a different domain/genre. However, social media text is diverse in domains and genre ranging from political to lifestyle discussions with short messages (e.g., tweets) and lengthy posts (e.g., blogs). Therefore, a system for analysing social media text needs to maintain consistent performance across domains/genres. This is a characteristic of the lexicon-based methods to sentiment classification. In this paper, we adopt the lexicon-based methods, hence, we concentrate on these methods in the rest of this related work section.

### 2.1. Lexicon-based methods

A lexicon-based sentiment analysis begins with the creation of a list of words associated with their sentiment polarity values (i.e. a sentiment lexicon), or the adoption of an existing one, from which the sentiment scores of terms are extracted and aggregated to predict sentiment of a given piece of text. Sentiment lexicons are either manually or semi-automatically generated from generic knowledge sources. Manually generated lexicons are obviously more accurate, however, they tend to have relatively low term coverage. In contrast, semi-automatically generated lexicons, such as by expanding a small set of seed words within a large corpus [6] or by dictionary propagation [4], have a high coverage of over 20,000 words. Moving away from traditional lexicons

that tend to capture individual terms, SenticNet has been introduced based on the idea of integrating concepts with common-sense knowledge [7]. SenticNet is a graph-structured resource with concepts as nodes and common-sense relationship between concepts as edges. Thus, when a concept extracted from a test text is triggered within SenticNet, common-sense knowledge associated with that concept can be exploited to enrich the machine's assessment of the problem being solved. Another resource with similar structure to SenticNet is WordNet [8], a machine readable dictionary that provides definitions of disambiguated word senses and establishes several relationships among them. These word senses were assigned quantified positive, negative and neutral polarity scores using an automated process to form the sentiment lexicon, SWN [4]. In this work, we use SWN as a general-purpose sentiment lexicon motivated by its relative high coverage of terms and its fine-grained sentiment information at word-sense level rather than term level.

A baseline lexicon-based classifier predicts the polarity class of a document using the aggregate of polarities of the terms contained in the document. With SWN, the sentiment dimension (positive or negative) that has the highest aggregate score becomes the sentiment class for the document [9–12]. This approach is inadequate for an effective sentiment analysis because the prior polarities of terms offered by a lexicon can be different from the contextual polarities of the terms. Such a difference, for instance, can arise due to the effect of linguistic rules such as negation or domain-specific term semantics that are not captured in a lexicon [13].

### 2.2. Contextual Analysis

This involves the adjustment of a term's prior polarity to reflect its polarity in a specific context. For example, the text *"I don't like the idea of smoking in general"* may be classified as positive because it is dominated by positive terms (*'like'* and *'idea'*). However, the appearance of the negation (*'don't'*) in the linguistic context of both terms rendered the text to be negative. In a contextual analysis strategy, the polarities of terms that are under the influence of negation are switched to the opposite sentiment dimension [14,15]. Similarly, polarity strength of terms that are under the influence of intensifiers (e.g., *'very', 'highly'*) or diminishers (e.g., *'slightly'* and *'a-little-bit'*) are increased and decreased respectively. Negation analysis is a particular challenge as the polarity of negated terms do not always translate to its opposite. For instance, whereas "It is *not good*" is more or less the same as "It is *bad*", "It is *not excellent*" is more positive than "It is *horrible*". Consequently, a shift approach was proposed as a preferred alternative to sentiment inversion for negation [16]. Here, the prior polarity of sentiment terms that are under the influence of negation is reduced by a certain weight, but the negation terms were not considered to bear sentiment of their own. However, a recent study suggests that negation terms are not just modifiers of sentiment but also indicators of sentiment [17]. In SWN, negation terms are associated with polarity scores. Thus, a strategy can be introduced to treat negation both as sentiment-bearing and as sentiment modifier for other terms.

Sentiment lexicons are typically generated independently of their target application. Thus, they tend to capture knowledge that is applicable across diverse domains (i.e. they are general-purpose). Not surprisingly deviations are common, especially on social media genres, due to variability in vocabulary usage resulting in poor sentiment coverage. Contextual deviations are also common, for instance where the sentiment polarities of terms differ from the domain-specific use of the terms. The poor sentiment coverage can be improved using a lexicon expansion strategy. In [18], a general-purpose lexicon has been expanded with Twitter-oriented sentiment-bearing terms extracted based on their mutual