



2014 Special Issue

A multi-label, semi-supervised classification approach applied to personality prediction in social media



Ana Carolina E.S. Lima, Leandro Nunes de Castro*

Natural Computing Laboratory, Mackenzie Presbyterian University, São Paulo, Brazil

ARTICLE INFO

Article history:

Available online 11 June 2014

Keywords:

Personality

Big Five

Social media

Twitter

Multi-label classification

Semi-supervised learning

ABSTRACT

Social media allow web users to create and share content pertaining to different subjects, exposing their activities, opinions, feelings and thoughts. In this context, online social media has attracted the interest of data scientists seeking to understand behaviours and trends, whilst collecting statistics for social sites. One potential application for these data is *personality prediction*, which aims to understand a user's behaviour within social media. Traditional personality prediction relies on users' profiles, their status updates, the messages they post, etc. Here, a personality prediction system for social media data is introduced that differs from most approaches in the literature, in that it works with groups of texts, instead of single texts, and does not take users' profiles into account. Also, the proposed approach extracts meta-attributes from texts and does not work directly with the content of the messages. The set of possible personality traits is taken from the Big Five model and allows the problem to be characterised as a multi-label classification task. The problem is then transformed into a set of five binary classification problems and solved by means of a semi-supervised learning approach, due to the difficulty in annotating the massive amounts of data generated in social media. In our implementation, the proposed system was trained with three well-known machine-learning algorithms, namely a Naïve Bayes classifier, a Support Vector Machine, and a Multilayer Perceptron neural network. The system was applied to predict the personality of Tweets taken from three datasets available in the literature, and resulted in an approximately 83% accurate prediction, with some of the personality traits presenting better individual classification rates than others.

© 2014 Elsevier Ltd. All rights reserved.

1. Introduction

Social media sites such as Twitter, Facebook and YouTube are based on human interaction and the concept of user-generated content. This leads to the creation and exchange of a vast amount of user-generated content, in a one-to-many type of communication, and entailing a massive production of free-form and interactive data (Barbier & Liu, 2011). Social media-oriented people tend to publish texts, audio or video pertinent to their lives that demonstrate their tastes and opinions. Statistics show that in the years 2013–2014, an average of 58 million Tweets were generated daily (Statistic Brain, 2014c); 3 million messages were sent every 20 min on Facebook (Statistic Brain, 2014a); 65 h of video were uploaded to YouTube every minute; 4 billion videos were viewed on YouTube (Statistic Brain, 2014d) every day; and 55 million photos were uploaded each day to Instagram (Statistic Brain, 2014b).

The data available within the social media platforms provide an unprecedented volume and richness of information about human behaviour and social interactions (Barbier and Liu, Data Mining in Social Media, 2011). Schultz and Schultz (2006) argue that the personality of an individual concerns his/her external characteristics and aspects that others can see. Social media provides data that makes it possible to understand who the users are and/or what they need. Thus, by analysing what is available in social media it is possible to identify important personality traits, that is, characteristics or qualities particular to a person, that can be used to describe his/her personality (Hall, Lindzey, & Campbell, 2000).

The Big Five (or Five Factor model) provides a parallel between the personality traits of an individual and his/her linguistic information, such as sentiment words. Social processes and family words can be used to define a computational model for predicting personality. This type of application can be seen in Golbeck, Robles, and Turner (2011b), Mairesse and Walker (2006), Mairesse, Walker, Mehl, and Moore (2007), Quercia, Kosinski, Stillwell, and Crowcroft (2011), Sumner, Byers, Boochever, and Park (2012) and Nunes, Teles, and Souza (2013), all of which rely on lexical

* Corresponding author. Tel.: +55 11 21148503.

E-mail addresses: aceslima@gmail.com (A.C.E.S. Lima), lnunes@mackenzie.br (L.N. de Castro).

analysis techniques to infer personality from Twitter and Facebook data. Lexical analysis can be supported by several studies that have attempted to understand the meaning of words carrying sentimental or emotional content (Kazemzadeh, Lee, & Narayanan, 2013; Tausczik & Pennebaker, 2010).

Building a model that can accurately predict personality using social media texts has several applications in a variety of areas such as marketing, business intelligence, psychology and sociology. It may also provide a qualitative perspective to the data, such as which sentiments are embedded in a message (Neri, Aliprandi, Capeci, Cuadros, & By, 2012) and therefore which products to recommend to the user (Dong, O'Mahony, Schaal, McCarthy, & Smyth, 2013; Feng & Qian, 2014). Data mining techniques together with concepts taken from psychology and applied to social media data can be extremely useful. For example, in understanding user behaviour (Adar, Weld, Bershad, & Gribble, 2007; Shi, Zhu, Cai, & Zhang, 2009; Xu, Zhang, Wu, & Yang, 2012); to study the dark triad (psychopathy, narcissism, and Machiavellianism) (Garcia & Sikström, 2013; Sumner et al., 2012; Wald, Khoshgoftaar, Napolitano, & Sumner, 2012); to identify criminal content (Lau, Xia, & Ye, 2014; Wang, Gerber, & Brown, 2012); to model affection (Martinez, Bengio, & Yannakakis, 2013); and to understand consumers' behaviour (Monteiro, Veiga, & Gonçalves, 2009).

From a computational point of view, the personality prediction based on the Big Five Model can be understood as a multi-label classification problem because an individual normally presents more than one personality trait, and each of these traits corresponds to a class for the classifier. Classification problems where there is no constraint on the number of labels that can be assigned to an object are called multi-label classification problems (Carvalho & Freitas, 2009). In the case of the Big Five Model personality is divided into five dimensions also known as the OCEAN model, which stands for Openness, Conscientiousness, Extroversion, Agreeableness and Neuroticism (Hall et al., 2000).

The present work introduces an approach based on the Big Five Model to predict personality in social media data, more specifically in groups of Twitter messages or 'Tweets'. The proposed approach is called PERSOMA (PERsonality prediction in SOcial Media datA). PERSOMA is composed of three main modules applied in cascade. In the first module, meta-data are extracted from Tweets to produce a 'meta-base'. No information about the users' profiles is required. Each of the five dimensions of the Big Five Model is seen as a personality trait, and, as each Tweet may contain from zero to five traits, the problem is characterised as a multi-label problem. In the second module, the multi-label problem is transformed into a set of five binary classification problems. Finally, as social media data are usually generated at a timescale much faster than can be annotated (i.e., labelled), the proposed approach makes use of a small number of labelled data and a semi-supervised learning approach to label (classify) the unlabelled data. In order to illustrate the use of PERSOMA, it is applied to a set of Tweets generated from a combination of three well-known sets from the social media data mining literature. Three different machine-learning algorithms are used as classifiers: Naïve Bayes (Feldman & Sanger, 2007); Support Vector Machine (SVM) (Hsu, Chang, & Lin, 2010); and a Multilayer Perceptron Neural Network (Haykin, 1998). The tool performs well and positive results are obtained that accurately predict the personality traits of 'Tweeters' without using their profile information, which is a different approach to that which appears in the literature (Golbeck, Robles, Edmondson, & Turner, 2011a; Quercia et al., 2011; Sumner et al., 2012).

The paper is organised as follows: Section 2 presents a discussion about personality and the Big Five Model; Section 3 provides a brief overview of the main efforts in the literature to perform personality prediction from social media data; Section 4 presents a detailed description of the proposed approach; Section 5

describes how the approach can be used in practice, by applying it to a dataset that combines three other datasets from the literature on social media mining; and Section 6 concludes the paper with a general discussion of the approach, its contributions and suggestions for future research.

2. On personality and the big five model

The word 'personality' comes from the Latin word *persona*, which refers to the mask used by actors in a theatre. This concept was derived from the understanding of personality as the combination of characteristics or qualities someone possesses. The first formal study of personality occurred within psychoanalysis, developed by Sigmund Freud. Freud divided the personality into three structures: the 'id', 'ego' and 'superego' in which id is the biological component of personality; ego is the rational component of the personality that acts according to the reality principle; and superego corresponds to the moral side of the personality, being composed of consciousness (Hall et al., 2000). Mairesse and Walker (2006) argue that personality can be described as a set of attributes that characterise an individual and involves behaviour, temperament, emotions, and the mind. However, this set of attributes can be relatively large. The diversity of attributes makes it difficult to evaluate personality because it does not provide a structure from which people can be classified and then compared. This same problem occurs when one tries to inform the sentiment present in a text (sentiment analysis) because the set of human emotions is large, thus making it difficult to select the basic emotions for a classification. To automate sentiment analysis, for instance, many researchers accept a simplified representation of sentiments according to their polarity (positive or negative) (Cambria, Schuller, Liu, Wang, & Havasi, 2013; Gangemi, Presutti, & Recupero, 2014).

A similar structure is observed in the study of personality. To make it possible to evaluate personality, various researchers have identified the most essential characteristics in order to create a personality model. Schultz and Schultz (2006) concluded that personality can be seen as a permanent and singular group of characteristics that can change in response to different situations. Thus, any personality prediction model must attempt to provide labels for these groups of characteristics. The first representation of a structure for personality was made in terms of traits, in which these represented a neuropsychic structure. Allport and Odbert (1936) defined two levels of traits, formulating a theory often referred to as the psychology traits (Hall et al., 2000). The first level is simply called *trait*, while the second is called *personal disposition* or *morphogenetic trait*. The traits, different from the personal disposition, are not peculiar to the individual, that is, a trait can be shared among various individuals. The study of traits allows us to make a comparison between individuals or groups of individuals, while in the study of personal disposition a single person is involved (Hall et al., 2000). Traits can be inferred from the frequency with which a person shows a certain type of behaviour, independent of the variety and intensity with which this behaviour appears (Hall et al., 2000). For example, it can be inferred that a person is sarcastic if he/she often makes sarcastic comments on a social media site.

A study by Allport and Odbert (1936) based on *lexical hypothesis* determined that the most important individual differences are encoded in language (Hall et al., 2000). According to the lexical hypothesis it is possible to systematise all behaviours and personality manifestations in terms (words), organised in a psychological dictionary of terms (Garcia, 2007). Cattell (1957) used the terms studied by Allport and Odbert in their analysis of the personality structure, and so began the definition of the Big Five Model or Five Factor Model. The Big Five describes a personality structure divided

Download English Version:

<https://daneshyari.com/en/article/404218>

Download Persian Version:

<https://daneshyari.com/article/404218>

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