



Estimation of the prevalence of adverse drug reactions from social media



Thin Nguyen^{a,*}, Mark E. Larsen^b, Bridianne O'Dea^b, Dinh Phung^a, Svetha Venkatesh^a, Helen Christensen^b

^a Centre for Pattern Recognition and Data Analytics, Deakin University, Australia

^b Black Dog Institute, University of New South Wales, Australia

ARTICLE INFO

Article history:

Received 27 August 2016

Received in revised form 23 February 2017

Accepted 21 March 2017

Keywords:

Consumer health informatics

Drug informatics

Adverse drug reactions

Social media

Word representation

Word embedding

ABSTRACT

This work aims to estimate the degree of adverse drug reactions (ADR) for psychiatric medications from social media, including Twitter, Reddit, and LiveJournal. Advances in lightning-fast cluster computing was employed to process large scale data, consisting of 6.4 terabytes of data containing 3.8 billion records from all the media. Rates of ADR were quantified using the SIDER database of drugs and side-effects, and an estimated ADR rate was based on the prevalence of discussion in the social media corpora. Agreement between these measures for a sample of ten popular psychiatric drugs was evaluated using the Pearson correlation coefficient, r , with values between 0.08 and 0.50. Word2vec, a novel neural learning framework, was utilized to improve the coverage of variants of ADR terms in the unstructured text by identifying syntactically or semantically similar terms. Improved correlation coefficients, between 0.29 and 0.59, demonstrates the capability of advanced techniques in machine learning to aid in the discovery of meaningful patterns from medical data, and social media data, at scale.

© 2017 Elsevier B.V. All rights reserved.

What was already known?

- Social media and other online generated content could help detecting adverse drug events (ADR).
- ADR lexicon could be built based on the National Library of Medicine's Medical Subject Heading (<http://www.nlm.nih.gov/mesh/>)
- The lexicon could be extended, either manually by looking up online or offline dictionaries, or automatically through using synonym packages or databases.

What this study has added?

- Confirm that social media could help estimating the prevalence of ADR efficiently.
- Word embedding techniques (word2vec) could help extending the lexicon of ADR terms automatically.

- The lexicon derived by word2vec improves the performance of using social media data to capture the prevalence of ADR.

1. Introduction

Advanced machine learning techniques, especially in natural language processing, have been employed to estimate the rate of adverse drug reactions (ADR) from social media [31]. Data of this type is often informal in nature, for example with lay or slang terms (for example, “*can't sleep*” or “*throwing up*”), rather than more formal terms (for example, “*insomnia*” or “*vomiting*”), which poses challenges for the task. To account for a variety of expressions of a single ADR, statistics on the co-occurrence with other terms might be examined. However, the dimension of the co-occurrence matrix increases with the size of the vocabulary, demanding increasing processing requirements. Thus, this approach becomes problematic when dealing with big data. Instead of capturing the global statistics from data, word2vec [28], a iteration-based framework, is suitable to suggest similar words at scale.

“*You shall know a word by the company it keeps*” (John Rupert Firth, 1957). Indeed, word2vec-like algorithms attempt to realize this principle, using deep learning approaches to capture relationships between words. In particular, word2vec is a neural network aiming to either predict surrounding words given a center word

* Corresponding author.

E-mail addresses: thin.nguyen@deakin.edu.au (T. Nguyen), mark.larsen@blackdog.org.au (M.E. Larsen), b.odea@blackdog.org.au (B. O'Dea), dinh.phung@deakin.edu.au (D. Phung), svetha.venkatesh@deakin.edu.au (S. Venkatesh), h.christensen@blackdog.org.au (H. Christensen).

(skip-gram model) or to predict a center word given its surrounding context words (Continuous Bag of Words model – CBOW) [28]. The resultant weights of the trained neural network become the representation for words in the vocabulary of the input corpus. In this representation, words are real valued vectors, namely word vectors. The representation is often called word embedding or distributed representation.

Word vectors are expected to reflect the relationships between words in the training corpus. Thus in the projected space, the spatial distance between the word vectors is related to the similarity in the context between the corresponding words. Specifically, the smaller the distance between word vectors, the closer the words are either in their syntax (for example, “apple” and “apples”) or semantics (for example, “clothing” and “shirt”).

Furthermore, using the means of vector arithmetic on the word vectors, one could answer analogous or relational questions, in form of “a is to b, as c is to ?”. Indeed, the word2vec framework has been shown to be the state of the art in capturing linguistic regularities, in both semantic and syntactic regularities [29]. For example, *bowl* is the answer to “clothing is to shirt as dish is to ?” as the word vector of *bowl* is closest to the word vector of (*clothing*–*shirt* + *dish*). Likewise, the word vector of *return* is closest to that of (*saw* – *see* + *returned*).

The first aim of this paper is to assess whether the rate of ADR described in social media documents is related to the known ADR rate. The second aim is to employ vector arithmetic on word vectors to identify additional terms for ADR, and to determine if this improves the rate of ADR detected. The current paper is organized as follows. Section 2 presents the background literature. Section 3 outlines the proposed methods, data, and experimental setup. Section 4 presents the results. Section 5 discusses these results and the limitations of the work, and Section 6 concludes the paper.

2. Background

2.1. Social media for health care

Social media is a core element of “Social Health” [4]. It has been integrated into medical practice and has reshaped health care services in several ways. This section highlights certain health care domains where social media has been utilized in.

2.1.1. Communication

Social media has improved health care quality with better communication between patients and clinicians. Through Facebook or Twitter, for example, social media provides a novel channel that quickly disseminates information to a large number of people virtually and at no cost [7]. For clinicians, there have been several physician-oriented social networking sites, such as Sermo¹. For patients, an example of a peer-to-peer health care system is PatientsLikeMe.² It offers patients an online platform to record and share health data about themselves through which people might learn efficient ways to deal with their own disease.

2.1.2. Health surveillance

Social media can be used to effectively build novel disease surveillance systems that detect, track and respond to infectious diseases, such as in the case of the 2009 H1N1 Influenza [6]. Social media could also help reach individuals with mental health disorders [22]. For example, postings referencing depression symptoms on Facebook were likely to report depression symptoms [30],

making Facebook a potential platform for major depressive disorder screening.

2.1.3. Promotional health

Promotional health is another area where social media can be employed. Effective promotion of healthy behaviors, such as weight loss programs [33], have been found to be feasible through social media. Social media is an inexpensive means to deliver behavioral change messaging and health promotion communications [20]. For example, Facebook and Twitter were found to be successful platforms to attract and engage a large number of users in sexual health promotion [34].

2.1.4. Medical intervention

Social media has been used for medical intervention. For example, Facebook was used to treat stress and depression in first-year medical students [12] and a cognitive-behavioral therapy delivered online incorporated into usual care can improve the treatment of depression [19].

2.1.5. Medical research

Social media has been utilized in health studies. For example, it has been leveraged to accelerate recruitment in clinical trials, leading to a new generation of large-scale clinical research [3]. It is proven to be feasible to conduct randomized, placebo-controlled, double-blinded trials entirely via the Internet, providing large pools of potential participants at reduced cost [17].

2.1.6. Clinical education

Clinicians have found that social media is a novel source of reference materials for medical education. For example, Facebook was used to deliver an HIV prevention program [7]; or YouTube has proven effective in many public health teaching programs [18].

2.1.7. Opinion mining for health care services

Social media has been employed for health ratings, aggregating millions of consumer reviews to find the best doctors and hospitals [14,13]. The online reviews and ratings of patients can be used to detect health care service quality, alleviating the cost of, or providing a supplemental source to, the traditional, standard survey approach [14].

2.1.8. Medical marketing

Social media has been used in medical marketing, constituting direct-to-consumer advertising which is a multi-billion-dollar business [25]. For example, social media was used for clinic branding and practice marketing for plastic surgeons [36].

2.2. Estimating adverse drug reactions from online texts

Social media and other online generated content has been employed for detecting of adverse drug events. For example, traces from online information-seeking could help detecting adverse events of drugs, suggesting a novel way for drug safety surveillance [35]. Other data sources, such as user-generated content in health forums or posts in health-related social networks, were also used to extract meaningful knowledge on patient safety [15,8,23]. Another approach is based on an online archive of scientific articles to automatically detect and validate adverse drug reactions. Examples include [5,32], where large-scale literature was mined to identify the adverse effects linked with prescription drugs. To improve the performance of adverse drug reaction detection, a fusion of multiple sources of information has been proposed. It could be a mix of chemical, biological and phenotypic information for drugs [26] or a combination of signals from official reports, such

¹ www.sermo.com.

² <http://www.patientslikeme.com>.

Download English Version:

<https://daneshyari.com/en/article/4966648>

Download Persian Version:

<https://daneshyari.com/article/4966648>

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