

Authorship Attribution in Health Forums

Victoria Bobicev

Technical University of Moldova
Chisinau, Moldova
vika@rol.md

Khaled El Emam

CHEO Research Institute, Ottawa
University of Ottawa, Ontario, Canada
kelemam@uottawa.ca

Marina Sokolova

CHEO Research Institute, Ottawa
University of Ottawa, Ontario, Canada
sokolova@uottawa.ca

Stan Matwin

Dalhousie University, Halifax, Canada
Polish Academy of Sciences, Warsaw, Poland
stan@cs.dal.ca

Abstract

The emergence of social media (networks, blogs, web forums) has given people numerous opportunities to share their personal stories, including details of their health. Although users mostly post under assumed nicknames, state-of-the-art text analysis techniques can combine texts from different media and use that linkage to identify private details of an individual's health. In this study we aim to empirically examine the accuracy of identifying authors of on-line posts on a medical forum.¹ Our results show a high accuracy of the authorship attribution, especially when text is represented by the orthographical features.

1 Introduction

Emergence of social media (networks, blogs, web forums) has given people numerous opportunities to share their personal stories, including details of their health (e.g., disease diagnosis, symptoms, treatment) (Velden and Emam, 2012; Bobicev et al, 2012):

- The transfer went well - my RE did it himself which was comforting. 2 embryos (grade 1 but slow in development) so I am not holding my breath for a positive.
- I've had 7 IUI and one ivf all cancelled due to not ovulating. I am a poor responder. What

bothers me the most is never getting to the point of actually going thru the procedure.²

Sharing personal health information (PHI) is a behavior that can be seen in 80% of Internet users, or in 59% of all adults, who reported searching for health information (Fox, 2011).

Although users mostly post under assumed nicknames, state-of-the-art text analysis techniques can combine texts from different forums and then use that linkage to identify private details of an individual's health. Aggregating and mining posts from five forums, Li et al. (2011) identified the user's full name, date of birth, spouse's name, home address, home phone number, cell phone number, email, occupation and the lab test results. The latter are highly indicative of the suspected disease, and hence, of the health conditions of the said individual.

In order to gauge how best to protect internet user anonymity, we first wanted to know the ability of Text Mining techniques in authorship attribution on medical forums, i.e. the task of identification of an author among other authors posting on the same forum. The attribution is based on comparison of a new text to texts previously written by known authors.

We obtained the empirical evidence on the posts from an on-line community of IVF (In Vitro Fertilization) patients. We achieved a highly accurate authorship attribution: up to 90% when the text is represented by the orthographical features.

¹ This work had been done when the first author was a visiting professor at CHEO Research Institute.

² The messages have an original spelling and punctuation.