# Detecting clinically related content in online patient posts

Courtland VanDam[a,*], Shaheen Kanthawala[a], Wanda Pratt[b], Joyce Chai[a], Jina Huh[c]

[a] Michigan State University, United States
[b] University of Washington, Seattle, United States
[c] University of California San Diego, United States

ABSTRACT

Patients with chronic health conditions use online health communities to seek support and information to help manage their condition. For clinically related topics, patients can benefit from getting opinions from clinical experts, and many are concerned about misinformation and biased information being spread online. However, a large volume of community posts makes it challenging for moderators and clinical experts, if there are any, to provide necessary information. Automatically identifying forum posts that need validated clinical resources can help online health communities efficiently manage content exchange. This automation can also assist patients in need of clinical expertise by getting proper help. We present our results on testing text classification models that efficiently and accurately identify community posts containing clinical topics. We annotated 1817 posts comprised of 4966 sentences of an existing online diabetes community. We found that our classifier performed the best (F-measure: 0.83, Precision: 0.79, Recall: 0.86) when using Naïve Bayes algorithm, unigrams, bigrams, trigrams, and MetaMap Symantic Types. Training took 5 s. The classification process took a fraction of 1 s. We applied our classifier to another online diabetes community, and the results were: F-measure: 0.63, Precision: 0.57, Recall: 0.71. Our results show our model is feasible to scale to other forums on identifying posts containing clinical topic with common errors properly addressed.

## 1. Introduction

Patients with chronic conditions visit online health communities to get help with managing their conditions [1]. In these communities, patients support one another through empathetic posts and consult on how to improve their daily health management strategies. At the same time, topics that can benefit from clinicians' expertise frequently appear in these patient discussions [2,3]. Messages containing such topics get buried in an overwhelming amount of posts, making it difficult for potential moderators to address them.

Moderators play an important role in online health communities. In addition to facilitating conversations, moderators add useful resources to posts containing clinically related questions [3]. Moderators also make sure information shared on their websites is not intended to be a substitute for professional medical advice by adding disclaimers or helping patients find relevant resources [3]. Patients self-moderate in online health communities where active, informal leaders exist [2,4]. For newly developing communities, however, such moderation activities around clinically related topics can be hard to do due to overwhelming amount of posts [5]. Efficiently identifying the patients' posts

needing additional, validated clinical resource would improve the quality of information shared in online health communities.

Many online health communities do not have moderators who can redirect questions to those with relevant expertise. Especially for those information needing clinical expertise, the community can benefit from knowing when certain questions need specific expertise over another. An automated system could be added by the forum owners to identify clinically-related posts to act upon it. If the information can be verified against a known knowledge base, e.g. WebMD, the system could respond to the user's post with either more information or additional verification that the advice is supported by the topic experts. If the information cannot be verified, or the concern can best be addressed by the user's physician, then the system could notify the user of the need that the content can benefit from extra verification as current moderators do [2].

In this paper, we develop a classification method to efficiently identify clinically related posts in online health communities. We examine specifically whether the clinical post addresses a medical question, a symptom, or a treatment. Existing work begins to address this problem, but the performance of classifiers could be improved [6]. The

classifier should also be able to scale to other communities. We used manually annotated data, feature design, feature selection methods, and comparisons across classifier algorithms to maximize the performance classifying clinically related posts in online diabetes communities. We also investigated the scalability of our classification model to other community context.

Our research questions include:

- How can feature design and selection techniques improve performance?
- Which classifier algorithm best perform in identifying topics from an online diabetes community?
- How high is the performance on detecting clinically related sentences in online health community posts?
- How much does our model built from one online health community generalize to other online health community contexts?

Below, we discuss related work, followed by the methods used to address these questions.

## 2. Related work

Online health communities present significant benefits to patients receiving support toward managing chronic disease. Research has shown the effects of using online social networks for chronic disease management. Merolli et al. summarized and analyzed the health outcomes and effects reported in previous studies [7]. One important benefit patients receive is support, both informational and emotional. Vlahovic et al. analyzed the satisfaction of users with their received support based on the type of support they requested, and they found that users seeking informational support and receiving emotional support were less satisfied than users seeking emotional support and receiving informational support [8]. De Choudhury et al. surveyed users about their sharing and seeking health-related information on Twitter [9]. They found 20% of the participants sought health-related information from Twitter. In particular, over half of those seeking information from Twitter were about seeking treatment information. Bui et al. found the sentiments of posts in online social support networks evolved from negative to positive sentiment [10]. Hartzler et al. investigated connecting patients based on their shared interests [11,12]. As such, existing work in online health social networks is focused on evaluating the efficacy of social support and devising ways to further augment support in online health communities. Further exploring work in improving qualities of sharing clinically related topics in online health communities can complement existing work around providing good quality social support to online health community members.

Huh et al. analyzed the roles of patients and moderators in online health communities [3]. They found that a majority of posts could benefit from clinical expertise, but there is not a sufficient number of clinical moderators to respond to all posts [13]. Even if moderators exist, sifting through a large number of community posts to identify posts needing clinical expertise can be overwhelming. To address this issue, a research team developed visualization tools to help moderators understand trends of aggregated online health community posts [14]. Furthermore, Huh et al. showed that moderators participate in online health communities to provide clinical expertise [3] and recommended patients to see a doctor [2]. A possible system to make these moderation activities more efficient is delivering moderators targeted posts needing their attention. To extract requirements for such system, Huh and Pratt interviewed clinicians while they read a subset of community threads to understand the challenges and necessary components of such a system [5]. The results indicated that clinicians identified clinically related keywords in posts as one of critical identifiers needing their attention and stated the importance of "triaging" the posts based on the severity of the problem expressed by the patients in their posts.

Researchers have attempted to identify clinically related posts in social media settings. McRoy et al. developed a classifier for community-based question answering websites, where the classification scheme included: factual clinical questions, patient-specific questions, and non-clinical questions [15]. Researchers also examined ways to identify authors of online health community posts-whether they are health professionals, which could inform the authority of clinical advice shared [16,17]. Abdaoui et al. used UMLS and other medical ontologies to determining whether the author of a post was a health professional or a lay man. Choumatare applied classification techniques to predict which patients had depression [18] to potentially provide help. Yang et al. used classification to detect posts that discuss adverse drug reactions [19]. Tuarob et al. classified whether or not each post from Twitter was health-related [20].Akbari et al. proposed an algorithm to detect wellness events, which are activities performed related to diet, exercise, or health [21].

As such, researchers have actively begun to investigate ways to deliver high quality information to patients online, augment social support, and provide interventions based on their stories posted online. Our work builds on this line of work, contributing new and improved ways to efficiently identify when patients need clinical expertise.

## 3. Methods

### 3.1. Data collection

Prior research has demonstrated that WebMD consists of active communities, where users discuss chronic health conditions [2,3,5,6]. WebMD is a health information portal website which provides information and tools to users for managing their health [22]. One critical feature of WebMD includes Exchanges, which is online communities where users discuss anything about managing their medical conditions. Each community is dedicated to one specific health condition, e.g. Diabetes or Heart Disease. We focused on the diabetes community (WDC) because it had the most active participation regarding balance between informational and emotional posts shared [6].

From WDC, we collected all threads posted between July 2007 (the beginning of the community) and July 2014 (the last date of data collection). A thread is a series of posts, which begins with a thread initiating post, followed by replies from other users. Because patients often initiate discussions in thread initiating posts [3], we examined only the posts that initiate threads of conversation through replies and replies to replies. We extracted 9576 thread initiating posts from the data we collected. We removed 538 duplicate posts. Each post contained one or more sentences. Fig. 1 demonstrates that most posts have 10 or fewer sentences. One post can consist of sentences that include clinically relevant keywords and those that do not. To simplify the scope what is considered a clinically relevant post, we designated each sentence as a unit of analysis. Our process is shown in Fig. 2.

We split all posts into sentences using Stanford's Natural Language Toolkit (NLTK) sentence tokenizer [23]. NLTK split posts into sentences by splitting on periods. We used a regular expression to identify and merge incorrectly split sentences into a whole sentence. Additionally, some users used other symbols, e.g. commas, to separate sentences. To address this, we manually identified and split these sentences during annotation.

We also collected posts from another online diabetes community (ODC2). This data was provided by our collaborator, who agreed to share their deidentified forum posts with us for the purpose of research and improving their own community. Identical to WDC, the community post structure was thread-based: each thread began with a thread initiating post, followed by the replies. We received 23,473 thread initiating posts from ODC2. We split these post into sentences using the same method described for WDC, which generated 2,009,005 sentences in total. To test the performance of our models, we applied our best performing classifier to all sentences from this data set. We then randomly selected 250 clinical sentences and 250 non-clinical sentences,