# Pharmaceutical Drugs Chatter on Online Social Networks

Matthew T. Wiley[a], Canghong Jin[b*], Vagelis Hristidis[a], Kevin M. Esterling[c]

**Correspondent:** Matthew Wiley, mwile001@cs.ucr.edu

[a]Department of Computer Science and Engineering, University of California, Riverside, CA, USA

[b]College of Computer Science and Technology, Zhejiang University, Hangzhou, China

[c]Department of Political Science, University of California, Riverside, CA, USA

*Work performed while visiting University of California, Riverside

**Abstract**

The ubiquity of Online Social Networks (OSNs) is creating new sources for healthcare information, particularly in the context of pharmaceutical drugs. We aimed to examine the impact of a given OSN's characteristics on the content of pharmaceutical drug discussions from that OSN. We compared the effect of four distinguishing characteristics from ten different OSNs on the content of their pharmaceutical drug discussions: (1) General vs. Health OSN; (2) OSN moderation; (3) OSN registration requirements; and (4) OSNs with a question and answer format. The effects of these characteristics were measured both quantitatively and qualitatively. Our results show that an OSN's characteristics indeed affect the content of its discussions. Based on their information needs, healthcare providers may use our findings to pick the right OSNs or to advise patients regarding their needs. Our results may also guide the creation of new and more effective domain-specific health OSNs. Further, future researchers of online healthcare content in OSNs may find our results informative while choosing OSNs as data sources. We reported several findings about the impact of OSN characteristics on the content of pharmaceutical drug discussion, and synthesized these findings into actionable items for both healthcare providers and future researchers of healthcare discussions on OSNs. Future research on the impact of OSN characteristics could include user demographics, quality and safety of information, and efficacy of OSN usage.

*Keywords:* Social Media, Health Social Media, Text Mining, Sentiment Analysis, Pharmaceutical Drugs, Frequent Itemsets.

## 1. Introduction

Numerous Online Social Networks (OSNs)[1] host Medicine 2.0 applications that focus specifically on user reviews of drugs [1-7]. Previous work has analyzed these discussions and confirmed that online drug reviews serve their purpose – i.e. users discuss medications and their effect on a disease or physical condition [8]. However, research is lacking on the impact of a given OSN's characteristics on the content of that OSN's discussions; e.g., if an OSN requires registration (e.g., providing an email address), does that affect the types of drugs users are willing to discuss?

Medicine 2.0 applications foster online communities where patients discuss their own healthcare decisions and experiences [9, 10]. These applications allow clinical researchers and citizen scientists to conduct crowdsourced health studies that complement traditional clinical trials in the public health research ecosystem [11, 12]. Such studies benefit other forms of knowledge generation, such as consumers' opinions of pharmaceutical drugs [13]. This knowledge is important: 24% of adults that use the Internet have read online reviews of a particular drug or medical treatment [14].

Moreover, there is increased interest from the research community in analyzing health-related content of OSNs. Previous work includes analyzing the content of health-related OSN discussions in terms of safety and quality, and detecting adverse drug reactions and events in OSN discussions; yet, previous work has not covered the impact of an OSN's characteristics on its discussions.

Therefore we analyzed the effect of four distinguishing characteristics of OSNs on a given OSN's content. These characteristics include: (1) OSN type – general (e.g. Twitter) versus health (e.g. WebMD); (2) if a given OSN moderates its posts; (3) if a given OSN requires registration; and (4) if a given OSN's discussions are in a Question and Answer (Q&A) format. We analyzed these characteristics both quantitatively (e.g., distribution of posts by drug type) and qualitatively (e.g., examining posts with the most frequent co-occurring medical concepts). Our results show that these OSN characteristics indeed affect the content of discussions related to pharmaceutical drugs. These effects include the type of discussions, the type of drugs discussed, the subjectivity of discussions, and the medical concept content.

In addition to the analysis results, this work also has the following key methodological contributions. We used sequences of carefully selected Web queries to identify important online drug review forums. We modified a

---

[1] **Abbreviation Note:** we use the term Online Social Networks (OSNs) to define social media platforms where users share content through messages; we further define these messages as posts. Examples of OSNs include Twitter and WebMD.

previous tool on medical concepts annotation to work on OSN posts. We enhanced the performance of an existing sentiment analysis dictionary to account for stemming and part of speech. We compared the drug distribution frequencies against a baseline, which assumes that all drugs have equal probability of being mentioned. Lastly, we mined OSN posts for frequent itemsets, where medical concepts were considered as items and each post is considered a transaction.

## 2. Related Work

Recently, there is increased interest in analyzing the content of health-related discussions in OSNs. Related work has chronicled the utility and potential benefit/harm of health-related discussions in OSNs; related work has focused on specific aspects of the information found in OSN discussions, but none focus on the impact of OSN characteristics. We demonstrate through our results that the characteristics of the OSNs adversely affect the type of content contained within each OSN. Coupling our findings with this related work provides possible (further) explanations of the findings from the related work. Another research area of recent interest at the intersection of healthcare and OSNs is detecting adverse drug events in OSN posts; the overreaching goal is real-time pharmacovigilance via the Internet. Our work complements this related work by giving further insight into the impact of OSN characteristics on discussions related to pharmaceutical drugs.

### 2.1 Analyzing Health Content of OSNs

Denecke and Nejdl [8] analyzed various Medicine 2.0 content and found that patient-authored postings contain more drug-related concepts than any other post. Further, they showed that drug reviews contain many disease related concepts and concluded that users searching for drugs or disorders will find results in patient-authored posts [8]. Lu *et al.* [15] studied the content of three discussion boards, from an online health community; they used one discussion board on diabetes and two on cancer. They found that drug-related postings accounted for a larger fraction of topics discussed on the diabetes board than the cancer boards [15].

Several works have looked at diabetes-related OSNs. Weitzman *et al.* [16] analyzed the quality and safety of diabetes-related OSNs and found that the quality/safety of information was variable across the ten sites under analysis. Shrank *et al.* [17] also qualitatively analyzed 15 diabetes-related OSNs – all of which feature a discussion or question forum – and they found a wide range in the number of members (from 3,000 to 300,000), one-third of the OSNs provided physicians answering questions, and two-thirds had site administrators reviewing posts. Zhang

*et al.* [18] analyzed posts from a Facebook diabetes group and found that over 60% of posts were providing information, followed by emotional support (17%) and eliciting information (12%).

Greene *et al.* [19] qualitatively analyzed the communications of Facebook communities dedicated to diabetes. They found many benefits for patients participating in these communities, such as community support and access to specialized knowledge, with little evidence of these communities supporting risky behaviors; however, one quarter of posts were explicit advertisements, some of which advertise non-FDA (Food and Drug Administration) approved products [19]. Two-thirds of posts were descriptions of personal experiences in diabetes management and a quarter of posts contained sensitive information unlikely to be revealed in doctor-patient interactions [19].

Goeuriot *et al.* [20] built and evaluated sentiment lexicons using drug reviews from a health social network. They built a general lexicon based on existing lexicons from the literature, and a domain lexicon based on drug reviews from the health social network. They showed that opinion mining of health social networks is possible, and using a combination of the general and domain lexicons achieves the best results [20].

*2.2 Detecting Adverse Events in OSNs*

Bian *et al.* [21] built two classifiers based on Twitter posts; one classifier to predict if a user (or someone they know) has used a particular drug, and a second classifier to classify if a post describes an adverse drug event. They obtain reasonable accuracy, but cite the noise in Twitter posts as one limitation to their approach [21]. Chee *et al.* [22] looked at predicting whether a drug will be withdrawn by the FDA using posts in Yahoo! Groups. While their classifier predicted many false positives (in the sense that a false positive is still on the market), a majority of the false positives with the greatest scores have been withdrawn from some market for a period of time [22].

Yang *et al.* [23] used association rule mining to detect adverse drug events in a health social network. Using data from the FDA, they confirmed correlations between drugs and adverse reactions in the posts [23]. Leaman *et al.* [24] validated that user comments from a health social network can be mined for adverse drug events. They built a lexicon based on manual annotations of users' posts and achieve reasonable accuracy using lexical matching [24].

**3. Methods**

*3.1 Datasets*

Our analysis used the ten OSNs listed in Table 1. Each of these OSNs was categorized as either a *general OSN* or a *health OSN*. General OSNs include Twitter, Google+, and Pinterest, which were chosen due to their popularity and various methods of sharing messages. To find health OSNs, we performed a series of Internet searches such as

"drug reviews", "user drug reviews", and "patient drug reviews"; these generic searches returned many results unrelated to drug reviews in social media, thus we used drug names from a list of the most popular drugs to find health OSNs, e.g., "Abilify reviews" and "Cymbalta reviews". We then chose the highest ranked sites that are public and have discussions by drug name. We only considered posts in health OSNs that originate from specific forums for reviewing drugs. Hence, posts from general forums or "Ask an Expert" forums were not collected from the health OSNs. Table A.1 of Appendix A lists the dates for which posts were collected and URLs for each OSN.

Each OSN was categorized further based on its moderation, registration requirements, and review format, as listed in Table 1; these categorizations are similar to related work that studies diabetes-related OSNs [16, 17]. This related work has shown that each of these categorizations is important: moderation affects the quality of information that is discussed; OSNs that require registration raises privacy concerns due to the poor readability of Privacy policies; and providing a forum where experts answer member questions best promotes safety for health OSNs. We consider an OSN to be moderated if a message is reviewed before becoming public. An OSN requires registration if it is necessary to create an account before publishing content. An OSN has a Q&A format if posts are formatted as comments/questions with replies/answers. We ignored categorizing each OSN based on whether users can posts anonymously, as this categorization is the same as the health versus general OSN category. Even if a health OSN requires registration, users have the option to post anonymously.

| Dataset | Health (H) or General (G)? | Moderated? | Registration Required? | Q&A Format? |
|---|---|---|---|---|
| Twitter | G | N | Y | N |
| Google+ | G | N | Y | N |
| Pinterest | G | N | Y | N |
| DailyStrength | H | N | Y | N |
| Drugs.com | H | Y | N | N |
| DrugLib.com | H | Y | N | N |
| everydayHealth | H | N | N | N |
| MediGuard | H | Y | Y | Y |
| medications | H | Y | Y | N |
| WebMD | H | N | N | N |

**Table 1 Various categorizations of each OSN. An OSN is moderated if a message is reviewed before becoming public. If registration is required, users must create an account before contributing content. An OSN is a Q&A format if reviews are formulated as comments/questions and replies/answers.**
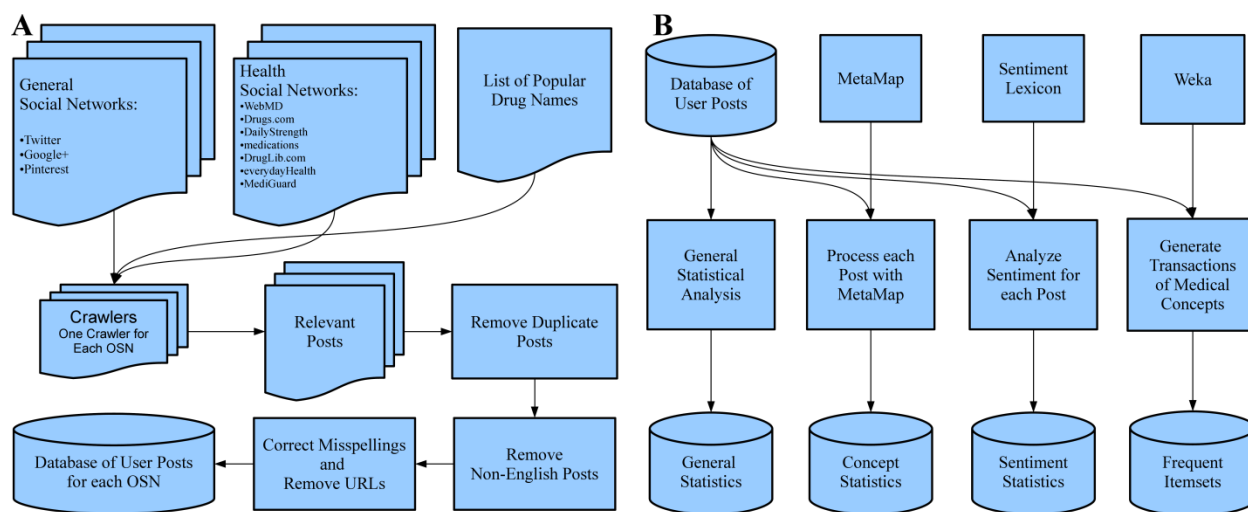
*3.2 Data Collection*

First we obtained a list of the 200 most popular drugs by prescriptions dispensed from RxList.com [25]. We then removed variants of the same drug (e.g., different milligram dosages) resulting in 122 unique drug names. This list was used as a filter for finding relevant posts. Posts from general OSNs were only considered relevant if one of the drug names was found in the post's text, whereas drug reviews from health OSNs were only collected for each of the 122 drugs. Note that most health OSNs will map equivalent drugs to the same drug review forum; for example, searching for Atorvastatin on DailyStrength will lead to the same series of drug reviews as Lipitor. The full list of drugs is given in Tables A.2, A.3, and A.4 of Appendix A.

For each OSN, we analyzed the layout of the website and built a crawler using Apache HttpComponents [26] – a library that enables web applications to obtain HTML content as if a web browser had downloaded and displayed the webpage; Twitter was handled separately using the Twitter API with the drug name list as a filter to collect matching tweets. Data for the rest of the OSNs was gathered by programmatically employing the search feature located on the respective OSN's website, where each drug name was specified as a query; e.g., we used Apache HttpComponents to search for Abilify on Google+. In the case of Pinterest and Google+, we collected all posts associated with the query; whereas the crawlers for health OSNs used the top search result that links to drug reviews (determining valid link patterns was done manually for each health OSN). The result is a series of HTML pages associated with a query for each OSN. Next, we extracted knowledge from each of the HTML pages using unique wrappers such as element id, location, or style. The wrappers and their content were extracted using jsoup, a Java HTML parser [27]. All pages for a given OSN follow the same HTML format, thus each of the wrappers were only defined once per OSN.

Posts in health social networks may contain metadata such as gender, age, length of membership, username, etc. However, even if a health OSN provides this information, the OSN allows users to leave this information blank; a manual inspection of posts on each of the health OSNs revealed that most users leave this information blank. Therefore we limited our data collection to the post text and date (if available). We collected all data in accordance with each OSN's terms of use, and therefore an OSN's data will not be made publicly available without first obtaining permission from the respective OSN.

Relevant posts obtained from the crawlers were further processed before the data analysis, as illustrated in Figure 1(A). First, duplicate posts are removed. Next, non-English posts are removed from the general OSNs (health OSNs only contained English posts); we used a Bayesian filter based on language profiles generated from

Wikipedia [28]. Next, we removed all hyperlinks and we corrected spelling mistakes in each of the posts; we corrected spelling errors using the first suggestion from HunSpell [29], an open source spell checker employed by several software packages. The result is a database of user posts that are relevant to the input list of prescription drug names for each OSN.



**Figure 1 (A) A visual overview of the data collection and preprocessing. Each crawler obtains a list of relevant posts using the OSNs as a seed and the list of drug names as a filter. These posts are then processed generating a database of English-only posts that have their spelling corrected. (B) An overview of the data analysis performed on the database of user posts. Four different types of results are generated by the data analysis: general statistics, concept statistics, sentiment statistics, and frequent itemsets.**

*3.3 Methods for Data Analysis*

The database created by the data collection process is then analyzed with four separate analyses: general statistics, medical concept statistics, sentiment statistics, and association rule mining. This process is illustrated in Figure 1(B). Since some OSNs have many more posts than others, we computed the average between each network when combining multiple OSNs into one result, rather than computing the average over all posts; otherwise, the results from Twitter or DailyStrength would decimate the results from each of the other OSNs.

*3.3.1 Methods for General Statistics*

One general statistic is the frequency of drugs based on their category. Drugs.com has a publicly available taxonomy of all drugs listed on its website [30], where one drug may be classified into multiple categories. We mapped our list of drug names to each of its top level categories as listed in the Drugs.com taxonomy; the distribution of these categories for our drug list is visualized in Figure A.1 of Appendix A. The full list of drug names along with their respective category or categories is given in Tables A.2, A.3, and A.4 of Appendix A.

For each OSN, we computed the frequency of each drug category and normalized this frequency by the total number of posts. For each OSN in a given category, we averaged the percentages of each drug category separately, and divided the sum of these percentages by the number of OSNs in the given category. Thus each OSN's distribution is weighted equally when presenting the distribution for the category. Otherwise, an OSN with many posts would dominate the category's distribution.

We analyzed OSN similarity by ranking the most frequent drugs. We measured similarity between each pair of ranked lists by using Spearman's footrule [31]. This measure of similarity considers the distance of each item (in terms of its rank) between two ranked lists. If the lists are identical, the value will be equal to zero, whereas a value of one denotes the maximum measure of disarray between the two lists. Other general statistics are presented in Appendix B.

*3.3.2 Methods for Medical Concept Statistics*

The MetaMap tool [32] was employed to annotate each post with medical concepts from the Unified Medical Language System (UMLS). The UMLS [33] is a compendium of several medical-focused ontologies. Thus MetaMap effectively represents each post as a set of medical concepts from the UMLS.

MetaMap was originally intended to annotate text for academic publications in the biomedical field, such as those available in PubMed. Related work has shown that MetaMap is not perfect for processing social media posts [34]. Thus, we manually inspected the annotations produced by MetaMap, and we removed annotations where MetaMap consistently misclassified UMLS concepts. A majority of mistakes were words that were misinterpreted as abbreviations in the social media posts. Other common mistakes included colloquial phrases not common to academic literature in the biomedical field. Some common mistakes include:

- the first-person narrative "I" was mapped to the UMLS concept for "Iodine" (C0021968)
- "so" was mapped to "Somalia" (C0021968)
- "fed" was mapped to "fish eye disease" (C0342895)

- "lol", "LOL" were mapped to "LOXL1 gene" (C1416898)

- "OMG" and "omg" were mapped to "OMG gene" (C1417949)

- "said" was mapped to "Simian Acquired Immunodeficiency Syndrome" (C0080151)

Mistakes similar to the ones given above were deleted from the MetaMap annotation results. We systematically analyzed each OSN by ordering every concept by its frequency and analyzing distinct phrases that were mapped for each concept. In total we identified 42 concepts that were incorrect. In general OSNs, these concepts accounted for over 5% of the total concept mappings, whereas these concepts account for less than 0.01% of the total concept mappings in health OSNs; the exact number of concept mappings (excluding mistakes) is reported in Appendix B.

Every concept in the UMLS is associated with one or more semantic types [35] (e.g., *Disease or Syndrome*). Each semantic type belongs to one of fifteen semantic groups [36], also defined by the UMLS. We analyzed the distribution of five semantic groups that relate to medical concepts, which include *Procedures*, *Disorders*, *Physiology*, *Chemicals and Drugs*, and *Anatomy*.

We considered the similarity of medical concept content between each OSN by ranking the most frequent semantic types. Again, we only considered semantic types that relate to medical concepts using the same five aforementioned semantic groups. We measured the similarity between each pair of ranked lists using Spearman's footrule; this is analogous to using Spearman's footrule for measuring OSN similarity with the most frequent drugs. Other medical concept statistics are presented in Appendix B.

### 3.3.3 Methods for Sentiment Statistics

The goal of sentiment analysis is to measure the average polarity and emotion of each post. Both are achieved by mapping phrases in each post to phrases from a sentiment lexicon. We use SentiWordNet [37], which contains a dictionary of phrases where each phrase is associated with a positive, negative, and objective score. Every term in SentiWordNet is subject to the constraint that the sum of the positive, negative, and objective score must equal one.

SentiWordNet distinguishes phrases based on their sense and part of speech. Therefore we tagged each word with its part of speech using the Stanford Core NLP tagger [38]. In order to remove variants of words, we stemmed both the posts and the terms in SentiWordNet; this was done to normalize words, e.g., rain, rains, and raining all become rain. Phrases form the posts are then mapped to phrases from SentiWordNet using the longest possible match first. In the case where one term has multiple senses, we averaged the score of all senses for the given term. We then computed the positive, negative, and objective scores of each post by averaging the scores from every

mapped term. The sentiment of a given OSN is measured by averaging the sentiment of all posts within that OSN. In the appendix we also present results from the NRC word-emotion lexicon [39] for analyzing the emotion of each OSN: negative–positive, anger–fear, trust–disgust, and anticipation–surprise.

*3.3.4 Methods for Frequent Itemsets*

Association rule mining is a data mining technique that learns relations between items given a database of transactions by first discovering frequent itemsets [40]. We applied this technique using UMLS concepts as items, where we considered each post to be a single transaction. Items were restricted based on their semantic groups; we analyzed frequent itemsets for medical concepts only and all UMLS concepts. Further, frequent itemsets were discovered separately for the health and general OSNs. For implementation we used the Weka machine learning toolkit [41]. Due to the large number of items and transactions, we employed the FP-growth algorithm [42] for discovering frequent itemsets. We removed trivial itemsets and only report itemsets that show interesting trends between categorizations of OSNs.
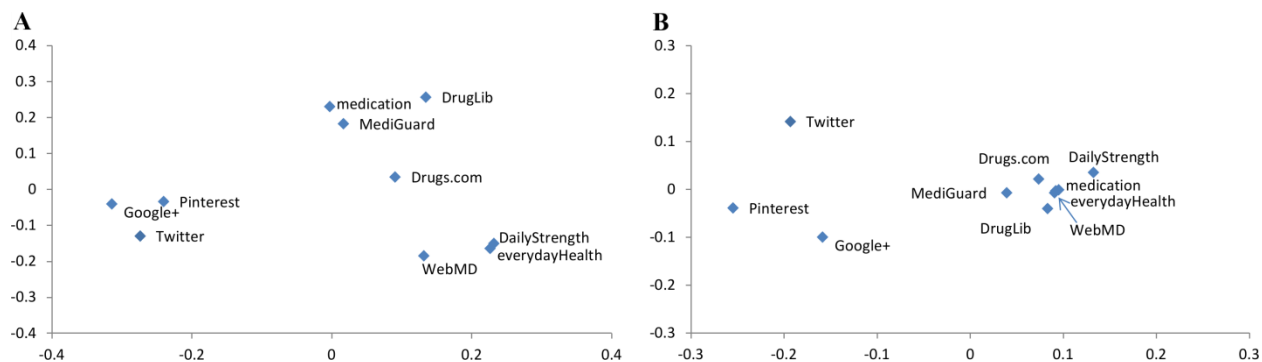
*3.3.5 Significance Testing*

For each of the aforementioned descriptive statistics, we conducted two statistical tests of significance. The first test we used is Pearson's Chi Squared Test for Independence [43]. The null hypothesis of this test is that there is no dependence between the variables in question, where the variables are the groupings of the OSNs, and thus the difference in distributions is due to random sampling. The alternative hypothesis is that there is some dependence between the groupings of OSNs. For each of the aforementioned statistics we built contingency tables and use the R programming language to compute the Chi statistic and p-value.

The second test we performed was Mann-Whitney U test, also known as the Wilcoxon-Mann-Whitney test [44]. For this test, we treat each post as an observation from the given grouping of OSNs. The null hypothesis is that the posts are drawn from the same population, whereas the alternative hypothesis is that one population tends to have larger values than the other. We also used the R programming language to compute the p-value for this test. Results of each test are reported with their corresponding figure, and detailed results of every test are given in Appendix G.

**4. Results**

Appendix B reports the statistics described in Section 3.3 for each OSN. Next, we compare the ten OSNs to each other using two measures of similarity. These measures include similarity between the most frequent drugs and the

most frequent semantic types using Spearman's footrule. The first measure shows which OSNs are similar based on the frequency of discussions about particular drugs, whereas the second measure shows which OSNs are similar based on the medical content (defined by the semantic types of the extracted concepts) in the discussions. Figure 2 illustrates these measures for each of the ten OSNs using metric multidimensional scaling [45].



**Figure 2 (A) multidimensional scaling of OSN similarity using Spearman's footrule with the top 25 most frequent drugs for each OSN. (B) multidimensional scaling of OSN similarity using Spearman's footrule with the top 30 semantic types for each OSN.**

As shown in Figure 2(A), there are three primary clusters of OSNs, with the general OSNs belonging to the bottom-left cluster, the non-moderated health OSNs belonging to the bottom-right cluster and the moderated health OSNs belonging to the top cluster. The reason for this clustering, also discussed in Sections 4.1 and 4.2, is that *these three groups mention different types of drugs*. The only OSN left out of these clusters is Drugs.com, which is a moderated health OSN; Drugs.com is separated from the other moderated health OSNs due to a higher number of psychotherapeutics in its top 25 drugs.

Figure 2(B) shows one cluster, which contains the health OSNs, and the three general OSNs separated from that cluster and each other. This figure suggests that *the medical content, in terms of UMLS semantic types, of health OSNs is similar, and differs from the medical content found in general OSNs*; further, this figure also suggests that the medical content in general OSNs varies across each OSN. For example, over 44%, 36%, and 45% of the concepts in Twitter, Google+, and Pinterest relate to Chemicals and Drugs respectively. Therefore Twitter is more likely to contain semantic types relating to Chemicals and Drugs in its top 25 semantic types.

The remainder of our results section examines each categorization of OSNs, and it is divided into four parts: (1) general versus health OSNs; (2) health OSNs that are non-moderated versus moderated; (3) health OSNs with registration versus no registration; and (4) health OSNs with a Q&A format versus health OSNs with a review

format. We omitted general OSNs from the last three categorizations of OSNs, since they all belong to the same categories (e.g., all are non-moderated).

*4.1 General versus Health OSNs*

Figure 3 compares the distributions of drug category frequency, polarity, and semantic groups of the health and general OSNs with the distribution of a uniform baseline. In Figure 3(A), this baseline is the distribution of the drug categories reported in Figure A.1. The baselines for Figures 3(B)-(C) assume a uniform distribution for all items matched in the database; e.g. the baseline in Figure 3(B) assumes a uniform distribution for all terms matched from SentiWordNet. All comparisons in this figure are significant with p < 0.001 for both significance tests.



**Figure 3 An overview of the analysis for general OSNs versus health OSNs: (A) the distribution of drug category frequencies; (B) the distribution of polarity; and (C) the distribution of semantic groups. Each baseline represents a uniform distribution: (A) assumes each drug from the drug list will appear with equal probability; (B) assumes each term mapped from SentiWordNet will appear with equal probability; and (C) assumes each UMLS concept extracted from the posts will appear with equal probability.**

Table 2 illustrates the major differences visualized in Figure 3. This table reports the highest absolute (i.e. ignoring sign) relative change of each item when compared to the baseline distributions. For example, there is a 590% increase in the number of posts related to genitourinary tract agents in general OSNs compared to the assumption that each drug would appear with equal probability. General OSNs have a decrease in both negative and positive polarity due to the number of objective terms in each post.

| Drug Category | | |
|---|---|---|
| Genitourinary Tract Agents | General | +590% |
| Nutritional Products | General | +290% |
| Psychotherapeutic Agents | Health | +167% |
| Nutritional Products | Health | -82% |
| Cardiovascular Agents | General | -74% |

| | | |
|---|---|---|
| Genitourinary Tract Agents | Health | -62% |
| Coagulation Modifiers | Health | -47% |
| Psychotherapeutic Agents | General | -39% |
| Coagulation Modifiers | General | +30% |
| Cardiovascular Agents | Health | -23% |
| **Polarity** | | |
| Negative | General | -32% |
| Positive | General | -18% |
| **Semantic Group** | | |
| Physiology | Health | +158% |
| Physiology | General | +60% |
| Chemical and Drugs | Health | -47% |
| Disorders | General | -30% |
| Chemical and Drugs | General | +23% |
| Disorders | Health | +12% |

**Table 2 Highest absolute relative changes of each item compared with the baselines shown in Figure 3. E.g., General Negative is computed as the difference between General Negative and Baseline Negative divided by Baseline Negative.**

Figure 3(A) shows some interesting trends between the types of drugs discussed in general and health OSNs. Firstly, both general and health OSNs have a smaller number of posts about cardiovascular agents compared to the baseline, and therefore *users of any OSN are less likely to post about cardiovascular agents such as Digoxin or Flomax*. The other drug categories show opposing trends between health and general OSNs – *drugs such as Viagra, Niaspan, and Warfarin are more common in general OSNs than drugs such as Cymbalta or Abilify, whereas the opposite is true for health OSNs*.

Figure 3(B) illustrates the differences in polarity between the health and general OSNs. *General OSNs use more objective terms; whereas health OSNs use more subjective terms*. There are several reasons for this result, and we are only able to speculate based on the data presented here. One possibility is that users of health OSNs are more likely to be serious patients who are suffering or recovering from serious problems. Another possibility is that the level of anonymity in health OSNs, where users often use name aliases, allows users to discuss more personal and subjective topics. Results for emotion, which are reported in Appendix C, show no significant differences between general and health OSNs.

Figure 3(C) illustrates the type of medical concepts discussed for general and health OSNs compared to a baseline that assumes each UMLS concept appears with equal probability. There is a large increase in the number

of concepts relating to physiology in health OSNs, but a decrease in the number of concepts relating to chemicals and drugs. General OSNs have more concepts relating to chemicals and drugs, and fewer concepts related to disorders. Further, these results suggest that *users of health OSNs are concerned with the effects of drugs on physiology, whereas users of general OSNs are either using drug names as slang or drug names in advertisements*.

*4.1.1 A Qualitative Analysis of General and Health OSNs*

Table 3 reports the most frequent itemsets of size 1 of medical concepts for health and general OSNs; itemsets of larger sizes are reported in Appendix C. *Health OSNs contain medical conditions, drug names and symptoms where the concept for sleep dominates* with a frequency of over 10%. *General OSNs contain many specific drugs names, where Viagra and Ibuprofen dominate* with frequencies over 27% and 16% respectively. Larger itemsets show that *general OSNs contain frequent itemsets of drugs that serve a similar purpose*; e.g., Ibuprofen, Tylenol, and Advil. *In general ONSs, drugs are often used as slang or in jokes*; e.g., "Viagra for women has been around for centuries. It's called money". *Funny news items are popular in general OSNs*; for example, Appendix C illustrates a series of frequent itemsets referring to Viagra, overdose, and amputated.

| Health OSNs | | General OSNs | |
|---|---|---|---|
| Sleep | 10.20% | Viagra | 27.20% |
| Depression | 4.81% | Ibuprofen | 16.89% |
| Headache | 4.11% | Penicillins | 4.61% |
| Tired | 4.02% | Sexual Intercourse | 2.36% |
| Weight Gain | 3.86% | Oxycodone | 2.26% |
| Anxiety | 3.62% | Sleep | 2.03% |
| Eating | 3.48% | Online Pharmaceutical Services | 1.66% |
| Mental Suffering | 3.22% | Acids | 1.50% |
| Dizziness | 3.17% | Headache | 1.45% |
| Lisinopril | 3.15% | Lactose Intolerance | 1.33% |

**Table 3 Frequent itemsets of size 1 for medical concepts.**

Table 4 reports frequent itemsets of size 1 of all concepts for health and general OSNs; itemsets of larger sizes are reported in Appendix C. Concepts for help, physician, milligram and started dominate health OSNs with frequencies greater than 12%, revealing that *users of health OSNs are discussing their experiences with their medications, and the differing strategies employed by their physicians*; e.g., "Because of my sleep troubles from Lexapro, [My doctor] started me on a new drug, Ambien to help me sleep with a dosage of 5 mg". *General ONSs contain posts from online pharmacies that advertise drugs for the best price with no prescription needed*; e.g., "[URL] with best price naprelan 250mg in internet rx overnight South Dakota". *Breaking news items about*

*pharmaceutical drugs are popular in general OSNs*; as discussed in Appendix C, the United States Food and Drug

Administration recommended lower dosages of Ambien for patients during a two week sample of Twitter data.
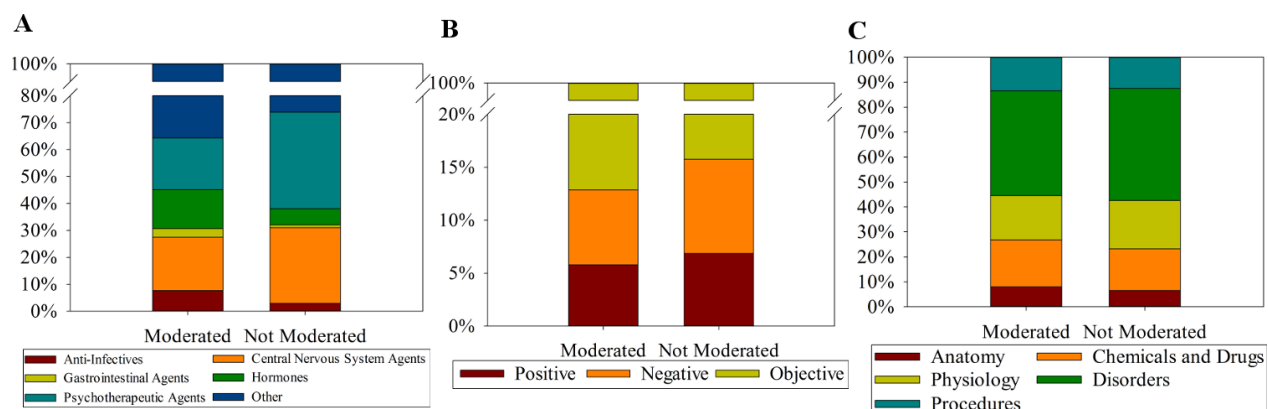
| Health OSNs | | General OSNs | |
|---|---|---|---|
| Help | 16.78% | Viagra | 22.42% |
| Physicians | 15.23% | Ibuprofen | 14.22% |
| Milligram | 13.75% | Milligram | 4.40% |
| Started | 12.24% | Penicillins | 3.83% |
| Sleep | 8.65% | Order | 3.15% |
| Dosage | 7.77% | Internet | 3.09% |
| Better | 7.57% | Prices | 2.23% |
| To be stopped | 5.50% | Overnight | 2.03% |
| Etiology aspects | 5.39% | Buying | 1.73% |
| Life | 5.33% | Sleep | 1.70% |

**Table 4 Frequent itemsets of size 1 for all UMLS concepts.**

*4.2 Moderated versus Non-moderated Health OSNs*

Figure 4 compares distributions of drug category frequency, polarity, and semantic groups of moderated and non-moderated health OSNs; all comparisons in this figure are significant with $p < 0.001$ for both significance tests.

Table 5 illustrates the major differences visualized in Figure 4. Appendix D reports the general statistics and medical concept statistics for moderated and non-moderated health OSNs.



**Figure 4 An overview of the analysis for moderated and not moderated OSNs. (A) The distribution of drug category frequencies; (B) the distribution of polarity; and (C) the distribution of semantic groups.**

| Drug Category | | |
|---|---|---|
| Gastrointestinal Agents | Moderated | +200% |
| Anti-infectives | Moderated | +158% |
| Respiratory Agents | Moderated | +143% |
| Hormones | Moderated | +141% |
| Psychotherapeutic Agents | Not Moderated | +87% |
| Central Nervous System Agents | Not Moderated | +41% |
| **Polarity** | | |
| Negative | Not Moderated | +25% |
| Positive | Not Moderated | +18% |

| Semantic Groups | | |
|---|---|---|
| Chemicals and Drugs | Moderated | +12% |

**Table 5 Highest absolute relative changes of each item for a given OSN group compared with the item of the other OSN grouping. E.g., Moderated Negative is computed as the difference between Moderated Negative and Not Moderated Negative divided by Not Moderated Negative.**

Figure 4(A) compares the distribution of drug categories between non-moderated health OSNs and moderated health OSNs. As noted in Table 5, *moderation affects the types of drugs users are willing to discuss; psychotherapeutic agents observed an 87% increase in frequency amongst non-moderated health OSNs.* Conversely, gastrointestinal agents, hormones, anti-infectives, and respiratory agents all observed an increase for moderated health OSNs, and a decrease for health OSNs that are not moderated.
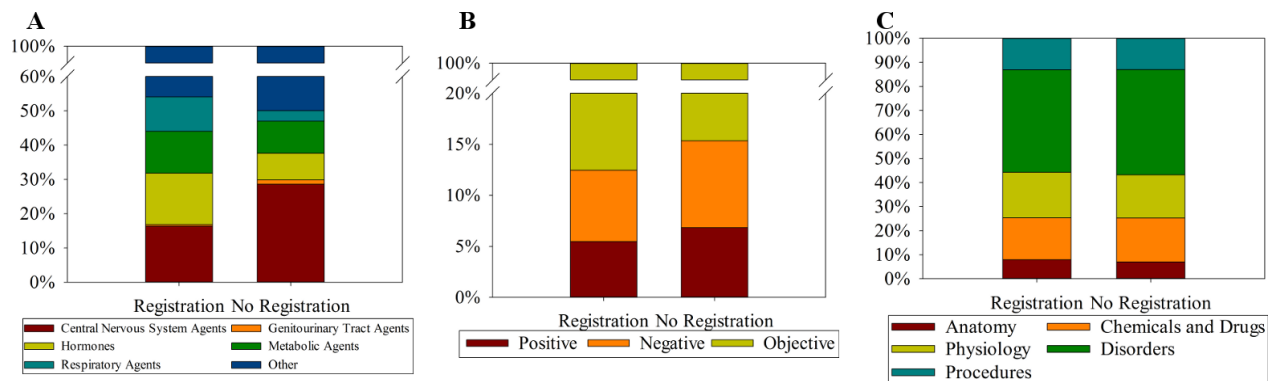
Figure 4(B) compares the distribution of polarity between health OSNs, non-moderated health OSNs, and moderated health OSNs. Also noted in Table 5, *moderation decreases the overall subjectivity, whereas non-moderated health OSNs increases subjectivity.* Thus, introducing moderation adds a level of objectivity to health OSNs.

Figure 4(C) reports the effect of moderation on semantic groups, and Appendix D reports the effect of moderation on emotion. *Overall, moderation has little effect on the medical concept content and emotional terms in health OSNs.* However, moderated health OSNs did have a slight increase on the number of terms relating to trust, whereas non-moderated health OSNs decreased the number of terms relating to trust. Further, moderated health OSNs increased the number of concepts relating to Chemicals and Drugs by 12. Appendix D reports frequent itemsets for health OSNs with and without moderation. These itemsets show that *users prefer non-moderated health OSNs when discussing psychotherapeutics and psychological conditions.*

*4.3 Registration versus no Registration in Health OSNs*

Figure 5 compares distributions of drug category frequency, polarity, and semantic groups of health OSNs that do or do not require registration; all comparisons in this figure are significant with $p < 0.001$ for both significance tests. Table 6 illustrates the major differences visualized in Figure 5. Appendix E reports the general statistics and medical concept statistics for health OSNs that do or do not require registration.

**Figure 5 An overview of the analysis for health OSNs that do or do not require registration. (A) The distribution of drug category frequencies; (B) the distribution of polarity; and (C) the distribution of semantic groups.**

| Drug Category | | |
|---|---|---|
| Respiratory Agents | Registration | +225% |
| Genitourinary Tract Agents | No Registration | +183% |
| Hormones | Registration | +92% |
| Central Nervous System Agents | No Registration | +74% |
| Metabolic Agents | Registration | +31% |
| **Polarity** | | |
| Positive | No Registration | +24% |
| Negative | No Registration | +21% |

**Table 6 Highest absolute relative changes of each item for a given OSN group compared with the item of the other OSN grouping. E.g., Registration Negative is computed as the difference between Registration Negative and No Registration Negative divided by No Registration Negative.**

Figure 5(A) compares the distribution of drug categories for health OSNs that do or do not require registration against all health OSNs as a baseline. As noted in Table 6, *registration affects the types of drugs users are willing to discuss; central nervous system agents observed a 74% increase in frequency amongst health OSNs that do not require registration.* Conversely, health OSNs that require registration have a 225% increase in posts about respiratory agents.

Figure 5(B) compares the distribution of polarity for health OSNs that do or do not require registration against all health OSNs as a baseline. *Similar to moderated health OSNs, requiring registration reduces the amount of subjectivity in health OSNs.*

Figure 5(C) reports the effect of registration on semantic groups, and Appendix E reports the effect of registration on emotion. *Overall, registration has little effect on the medical concept content and emotional terms in health OSNs.* Appendix E reports frequent itemsets for health OSNs that do or do not require registration. Similar

to moderation, these itemsets show that *users prefer health OSNs that do not require registration when discussing psychotherapeutics and psychological conditions*.

*4.4 Review versus Q&A format*

Figure 6 compares distributions of drug category frequency, polarity, and semantic groups of health OSNs that have a review format with health OSNs that have a Q&A format; all comparisons in this figure are significant with p < 0.001 for both significance tests. Table 7 illustrates the major differences visualized in Figure 6. Appendix F reports the general statistics and medical concept statistics for health OSNs with a review or Q&A format.



**Figure 6 An overview of the analysis for health OSNs with a review or Q&A format. (A) The distribution of drug category frequencies; (B) the distribution of polarity; (C) the distribution of semantic groups.**

| Drug Category | | |
|---|---|---|
| Anti-infectives | Review | +354% |
| Coagulation Modifiers | Q&A | +243% |
| Metabolic Agents | Q&A | +63% |
| Gastrointestinal Agents | Review | -52% |
| Psychotherapeutic Agents | Review | +47% |
| **Polarity** | | |
| Negative | Review | +144% |
| Positive | Review | +110% |
| **Semantic Groups** | | |
| Chemicals and Drugs | Q&A | +36% |
| Procedures | Q&A | +24% |
| Disorders | Review | +18% |
| Physiology | Review | +16% |

**Table 7 Highest absolute relative changes of each item for a given**

**OSN group compared with the item of the other OSN grouping. E.g.,**

**Review Negative is computed as the difference between Review**

**Negative and Q&A Negative divided by Q&A Negative.**

Figure 6(A) compares the distribution of drug categories for health OSNs that have a review format or Q&A format. Health OSNs that have a Q&A format have a 243% and 63% increase in posts related to coagulation modifiers and metabolic agents respectively. Posts about psychotherapeutic agents and anti-infectives observed an increase of 47% and 354% in health OSNs with a review format. This suggests that *users are less likely to ask questions about Abilify or Penicillin, but users are more likely to ask questions about Warfarin, Advair, or Lipitor*.

Figure 6(B) compares the distribution of polarity for health OSNs that have a review format or Q&A format. Health OSNs with a Q&A format are much more objective than health OSNs with a review format, where health OSNs with a review format observed an increase of 144% and 110% to negativity and positivity respectively. Thus, *users of health OSNs with a Q&A format tend to post in an objective manner, rather than subjective opinions regarding a particular drug*.

Figure 6(C) compares the distribution of semantic groups for health OSNs that have a review format or Q&A format. Health OSNs with a Q&A format observed an increase of 36% and 24% for Chemicals and Drugs and Procedures respectively; whereas health OSNs with a review format observed an increase of 18% and 16% to Disorders and Physiology respectively. This suggests *users ask questions that focus on drugs and procedures rather than questions about specific disorders or effects on their physiology*.

**5. Discussion**

Our results section has demonstrated the similarities and differences of OSNs in the context of pharmaceutical chatter in OSNs. Together, these data may help inform patients and healthcare providers about the type of content related to pharmaceutical drugs on OSNs. As pointed out by Eysenbach, OSNs (including health OSNs) are essentially an apomediated environment [10], where users take over the role of intermediary and guide other users to relevant and accurate information.

Based on our findings, healthcare providers could advise patients on the use of OSNs. Examples include: the prevalence and legitimacy of online pharmacies due to the high number of advertisements from online pharmacies in general OSNs; general OSNs are good sources of breaking news, particularly if that news was reported by a trusted source such as United States Food and Drug Administration; thousands of other patients are discussing health conditions and their treatments on health OSNs, yet these discussions may be subjective or biased; health OSNs that require registration, have moderation, or a Q&A format tend to be more objective, and thus information is less opinionated.

Our results may also guide the creation of new and more effective domain-specific health OSNs. Furthermore, these data may help future researchers that study OSNs make informed decisions about the social networks chosen for study when consider health content in OSNs. In the context of pharmaceutical drug chatter in OSNs: general OSNs are sources of jokes, news, and advertisements; health OSNs are sources of user experiences' with pharmaceutical drugs and strategies employed by their physicians for a particular medical condition or set of medical conditions; also, sleep and sleep related problems are a common theme throughout health OSNs. Drugs and diseases relating to the brain or central nervous system are more frequently discussed on health OSNs that are non-moderation and do not require registration respectively. In contrast, more prevalent diseases, such as asthma, hypertension, or high cholesterol are more frequently discussed on health OSNs that have moderation or require registration. Lastly, users are more likely to ask questions in public spaces about respiratory agents and hormones.

*5.1 Limitations*

We did not consider demographics of users in this study as this information was not present in every source. Therefore we cannot generalize our results to the general population. However, given that nearly 1 in 4 adults in 2011 that used the Internet, also looked for reviews on drugs or medical treatments [14], we argue that our results are still consequential to a substantial portion of the general population.

Another limitation of our work is that we did not remove messages that would be considered spam. The definition of spam is subjective – health social networks would remove pharmaceutical advertisements, whereas general social networks would not remove these advertisements from verifiable companies. We manually examined over 1,000 posts from health OSNs, and there was no evidence of any advertisements or spam in these OSNs. Moderated health OSNs would prevent messages from being published if a message was an advertisement or spam. Health ONS that are not moderated contain features for users to report messages as spam; thus these messages would be removed at some point after their publication.

General OSNs take steps to eliminate spam [46-48], but these OSNs clearly contain pharmaceutical advertisements. We believe it is worthwhile to consider these advertisements when examining general OSNs, as any user (or researcher) may be exposed to posts advertising overnight prescriptions for controlled substances. Further, we assert that including advertisements do not materially affect our results, for several reasons. First, our frequent itemset analysis revealed that tweets containing drug names from advertisements (e.g., Viagra or Ibuprofen), are also contained in tweets from real users. Second, Twitter restricts its policy for advertising of health and pharmaceutical

products [49], and Twitter's policy on ads specifically states that ads for illegal goods and services are prohibited [50]. And third, manual examination of Google+ and Pinterest found that these datasets contain far fewer pharmaceutical advertisements than Twitter. For all of these reasons, we have chosen not to exclude advertisements from our data. In our future work, we plan to build an advertisement classifier to study the role of advertisements in health-related OSN chatter.

There are also technical limitations with our approach. Due to the volume of Twitter posts, we only selected a ten month sample of posts, whereas we collected as many posts as possible for each of the other datasets. Ideally, we would examine all posts from Twitter since Twitter's beginning. Due to crawling constraints, we did not consider every social network where users post messages with respect to pharmaceutical drugs. MetaMap is not perfect for annotating social media posts, but we did clean up its output by removing annotations that are obviously incorrect. While the UMLS is a compendium of several medically focused ontologies, an ideal ontology for OSN posts about pharmaceutical drugs would be built using a specialized lexicon for health-related posts in social media; such a lexicon would also apply to the sentiment lexicons, where terms such as "omg" and "lol" are not mapped to any word in each of the sentiment lexicons used in this work.

## 6. Conclusion

With the objective to analyze the impact of OSN characteristics on the content of pharmaceutical drug discussions, we have reported several patterns of information from ten different OSNs. We demonstrated that an OSN's characteristics affect the type of discussions, the type of drugs discussed, the subjectivity of discussions, and the medical concept content. We synthesized these findings and proposed actionable items for both healthcare providers and future researchers of healthcare discussions on OSNs. Future research on the effect of OSN characteristics in healthcare discussions could include user demographics, quality and safety of information, and efficacy of OSN usage.

**References**

[1] DailyStrength. http://dailystrength.org (Accessed: January 31, 2013).
[2] DrugLib.com. http://www.druglib.com (Accessed: January 31, 2013).
[3] Drugs.com. http://www.drugs.com (Accessed: January 31, 2013).
[4] everydayHealth. http://www.everydayhealth.com (Accessed: January 31, 2013).
[5] medications.com. http://medications.com (Accessed: January 31, 2013).
[6] MediGuard. http://www.mediguard.org (Accessed: January 31, 2013).
[7] WebMD. http://www.webmd.com/ (Accessed: January 31, 2013).
[8] Denecke K, Nejdl W. How valuable is medical social media data? Content analysis of the medical web. J Inf Sci. 2009;179:1870-80. http://dx.doi.org/10.1016/j.ins.2009.01.025
[9] Van De Belt TH, Engelen LJ, Berben SA, Schoonhoven L. Definition of Health 2.0 and Medicine 2.0: a systematic review. J Med Interenet Res. 2010;12:e18. http://dx.doi.org/10.2196/jmir.1350
[10] Eysenbach G. Medicine 2.0: social networking, collaboration, participation, apomediation, and openness. J Med Internet Res. 2008;10:e22. http://dx.doi.org/10.2196/jmir.1030
[11] Swan M. Scaling crowdsourced health studies: the emergence of a new form of contract research organization. J Pers Med. 2012;9:223-34. http://dx.doi.org/10.2217/pme.11.97
[12] Swan M. Crowdsourced health research studies: an important emerging complement to clinical trials in the public health research ecosystem. J Med Internet Res. 2012;14:e46. http://dx.doi.org/10.2196/jmir.1988
[13] Swan M. Health 2050: The Realization of Personalized Medicine through Crowdsourcing, the Quantified Self, and the Participatory Biocitizen. J Pers Med. 2012;2:93-118. http://dx.doi.org/10.3390/jpm2030093
[14] Fox S. The social life of health information, 2011. Pew Internet & American Life Project; 2011.
[15] Lu Y, Zhang P, Liu J, Li J, Deng S. Health-Related Hot Topic Detection in Online Communities Using Text Clustering. PloS One. 2013;8:e56221. http://dx.doi.org/10.1371/journal.pone.0056221
[16] Weitzman ER, Cole E, Kaci L, Mandl KD. Social but safe? Quality and safety of diabetes-related online social networks. J Am Med Inform Assoc. 2011;18:292-7. http://dx.doi.org/10.1136/jamia.2010.009712
[17] Shrank WH, Choudhry NK, Swanton K, Jain S, Greene JA, Harlam B, et al. Variations in structure and content of online social networks for patients with diabetes. Arch Intern Med. 2011;171:1589. http://dx.doi.org/10.1001/archinternmed.2011.407
[18] Zhang Y, He D, Sang Y. Facebook as a Platform for Health Information and Communication: A Case Study of a Diabetes Group. J Med Syst. 2013;37:1-12. http://dx.doi.org/10.1007/s10916-013-9942-7
[19] Greene JA, Choudhry NK, Kilabuk E, Shrank WH. Online social networking by patients with diabetes: a qualitative evaluation of communication with Facebook. J Gen Intern Med. 2011;26:287-92. http://dx.doi.org/10.1007/s11606-010-1526-3
[20] Goeuriot L, Na J-C, Min Kyaing WY, Khoo C, Chang Y-K, Theng Y-L, et al. Sentiment lexicons for health-related opinion mining. Proceedings of the 2nd ACM SIGHIT International Health Informatics Symposium. Miami, FL: ACM; 2012. p. 219-26. http://dx.doi.org/10.1145/2110363.2110390
[21] Bian J, Topaloglu U, Yu F. Towards large-scale twitter mining for drug-related adverse events. Proceedings of the 2012 International Workshop on Smart Health and Wellbeing. Maui, Hawaii: ACM; 2012. p. 25-32. http://dx.doi.org/10.1145/2389707.2389713
[22] Chee BW, Berlin R, Schatz B. Predicting adverse drug events from personal health messages. AMIA Annual Symposium Proceedings. Washington D.C.: American Medical Informatics Association; 2011. p. 217-26.
[23] Yang CC, Jiang L, Yang H, Tang X. Detecting Signals of Adverse Drug Reactions from Health Consumer Contributed Content in Social Media. Proceedings of ACM SIGKDD Workshop on Health Informatics. Beijing, China: ACM; 2012.
[24] Leaman R, Wojtulewicz L, Sullivan R, Skariah A, Yang J, Gonzalez G. Towards internet-age pharmacovigilance: extracting adverse drug reactions from user posts to health-related social networks. Proceedings of the 2010 Workshop on Biomedical Natural Language Processing. Uppsala, Sweeden: Association for Computational Linguistics; 2010. p. 117-25.
[25] RxList. http://www.rxlist.com/script/main/hp.asp (Accessed: January 15, 2013).
[26] Apache. Apache HttpComponents. http://hc.apache.org/ (Accessed: January 4, 2013).
[27] Hedley J. jsoup: Java html parser. http://jsoup.org/ (Accessed: January 4, 2013).
[28] Shuyo N. language-detection - Language Detection Library for Java. http://code.google.com/p/language-detection/ (Accessed: February 25, 2013).
[29] Németh L. Hunspell. http://hunspell.sourceforge.net/ (Accessed: 25 Feb, 2013).
[30] Drugs.com. Drugs By Category. http://www.drugs.com/drug-classes.html?tree=1 (Accessed: March 4, 2013).
[31] Fagin R, Kumar R, Sivakumar D. Comparing top k lists. SIAM J on Discret Math. 2003;17:134-60.

[32] Aronson AR. Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program. AMIA Annual Symposium Proceedings. Washington D.C.: American Medical Informatics Association; 2001. p. 17-21.

[33] Bodenreider O. The unified medical language system (UMLS): integrating biomedical terminology. Nucleic Acids Res. 2004;32:D267-D70. http://dx.doi.org/10.1093/nar/gkh061

[34] Denecke K, Soltani N. The Burgeoning of Medical Social-Media Postings and the Need for Improved Natural Language Mapping Tools. Where Humans Meet Machines: Springer; 2013. p. 27-43.

[35] National Library of Medicine. Current Semantic Types. http://www.nlm.nih.gov/research/umls/META3_current_semantic_types.html (Accessed: March 12, 2013).

[36] National Library of Medicine. The UMLS Semantic Groups. http://semanticnetwork.nlm.nih.gov/SemGroups/ (Accessed: April 2, 2013).

[37] Baccianella S, Esuli A, Sebastiani F. Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10). Valletta, Malta: European Language Resources Association (ELRA); 2010.

[38] Toutanova K, Klein D, Manning CD, Singer Y. Feature-rich part-of-speech tagging with a cyclic dependency network. Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1. Edmonton, Canada: Association for Computational Linguistics; 2003. p. 173-80. http://dx.doi.org/10.3115/1073445.1073478

[39] Mohammad SM, Turney PD. Crowdsourcing a word–emotion association lexicon. Comput Intell. 2012. http://dx.doi.org/10.1111/j.1467-8640.2012.00460.x

[40] Agrawal R, Imieliński T, Swami A. Mining association rules between sets of items in large databases. ACM SIGMOD Record: ACM; 1993. p. 207-16. http://dx.doi.org/10.1145/170036.170072

[41] Hall M, Frank E, Holmes G, Pfahringer B, Reutemann P, Witten IH. The WEKA data mining software: an update. ACM SIGKDD Explorations Newsletter. 2009;11:10-8. http://dx.doi.org/10.1145/1656274.1656278

[42] Han J, Pei J, Yin Y. Mining frequent patterns without candidate generation. ACM SIGMOD Record: ACM; 2000. p. 1-12. http://dx.doi.org/10.1145/335191.335372

[43] Pearson K. X. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science. 1900;50:157-75.

[44] Fay MP, Proschan MA. Wilcoxon-Mann-Whitney or t-test? On assumptions for hypothesis tests and multiple interpretations of decision rules. Statistics surveys. 2010;4:1.

[45] Davison ML. Multidimensional scaling. New York: Wiley; 1983.

[46] Twitter Blog. Shutting Down Spammers. http://blog.twitter.com/2012/04/shutting-down-spammers.html (Accessed: April 2, 2013).

[47] Pinterest Blog. Fighting Spam. http://blog.pinterest.com/post/37347668045/fighting-spam (Accessed: April 2, 2013).

[48] Chansanchai A. Google+ wages war on spam comments. http://www.nbcnews.com/technology/technolog/google-wages-war-spam-comments-277331 (Accessed: April 2, 2013).

[49] Twitter. Health and pharmaceuitcal products and services. http://support.twitter.com/articles/20170441-health-and-pharmaceutical-products-and-services (Accessed: October 15, 2013).

[50] Tiwtter. Twitter Ads Policies. http://support.twitter.com/groups/56-policies-violations/topics/236-twitter-rules-policies/articles/20169693-twitter-ads-policies (Accessed: October 15, 2013).

# Appendix A  Online Social Network and Drug Summary

*A.1 Online Social Network Summary*

Table A.1 lists each of the ten Online Social Networks (OSNs) investigated in this work with their respective

website, and the start and end dates of posts collected from each OSN.  Not every OSN marks posts with

timestamps, therefore these networks are marked with the date they were crawled. For Twitter, we use a random

sample of 20% of the posted tweets for text analysis purposes, due to its large volume, but we use all tweets to

report the average number of posts per day in Table B.1.

| Dataset | URL | Start | End |
|---|---|---|---|
| Twitter | www.twitter.com | Dec. 29, 2012 | Oct. 31, 2013 |
| Google+ | plus.google.com | Jan. 1, 2011 | Jan. 31, 2013 |
| Pinterest | www.pinterest.com | N/A | Feb. 11, 2013* |
| DailyStrength | www.dailystrength.org | N/A | Jan. 15, 2013* |
| Drugs.com | www.drugs.com | Apr. 2, 2007 | Jan. 23, 2013 |
| DrugLib.com | www.druglib.com | N/A | Feb. 11, 2013* |
| everydayHealth | www.everydayhealth.com | Jan. 2, 2001 | Jan. 31, 2013 |
| MediGuard | www.mediguard.org | Jan. 21, 2007 | Jan. 31, 2013 |
| medications | www.medications.com | N/A | Feb. 13, 2013* |
| WebMD | www.webmd.com | Sept. 18, 2007 | Jan. 19, 2013 |

**Table A.1  An overview of the OSNs analyzed in this work.  The start and end dates**

**represent the timestamp of the first and last post from each dataset.**

**\* The date an OSN was crawled for OSNs that do not mark posts with an exact timestamp.**

*A.2 Drug Summary*

Tables A.2, A.3 and A.4 list the most popular drugs by prescriptions dispensed, as given on RxList.com [1].  Each of

these drugs was classified into one or more drug groups, according to the drug taxonomy available on Drugs.com

[2].  Each drug is associated with one or more categories.

| Gastrointestinal Agents | Genitourinary Tract Agents | Topical Agents | Alternative Medicines | Nutritional Products | Coagulation Modifiers |
|---|---|---|---|---|---|
| Famotidine | Cialis | Mupirocin | Lovaza | Folic | Plavix |
| Nexium | Detrol | Nasonex | | Klor-Con | Warfarin |
| Omeprazole | Viagra | Premarin | | Niaspan | |
| Pantoprazole | | Xalatan | | | |
| Ranitidine | | | | | |

**Table A.2 Listing of drugs that were classified as Gastrointestinal Agents, Genitourinary Tract Agents,**

**Topical Agents, Alternative Medicines, Nutritional Products, and Coagulation Modifiers.**

| Hormones | Anti-infectives | Psychotherapeutic Agents | Respiratory Agents |
|---|---|---|---|
| Levothyroxine | Amoxicillin | Abilify | Advair |

| Levoxyl | Azithromycin | Amitriptyline | Albuterol |
|---|---|---|---|
| Loestrin | Cefdinir | Citalopram | Cheratussin |
| Methylprednisolone | Cephalexin | Cymbalta | Combivent |
| NuvaRing | Ciprofloxacin | Effexor | Fexofenadine |
| Ocella | Doxycycline | Fluoxetine | Flovent |
| Prednisone | Fluconazole | Lexapro | Fluticasone |
| Premarin | Levaquin | Paroxetine | Hydrocodone |
| Synthroid | Penicillin | Seroquel | Proair |
| TriNessa | Sulfamethoxazole | Sertraline | Promethazine |
| | | Trazodone | Proventil |
| | | Zyprexa | Singulair |
| | | | Spiriva |
| | | | Ventolin |

**Table A.3 Listing of drugs that were classified as Hormones, Anti-infectives, Psychotherapeutic Agents, and**

**Respiratory Agents.**

| Metabolic Agents | Cardiovascular Agents | Central Nervous System Agents |
|---|---|---|
| Actonel | Amlodipine | Alprazolam |
| Actos | Atenolol | Ambien |
| Alendronate | Benazepril | Amphetamine |
| Allopurinol | Benicar | Aricept |
| Crestor | Carvedilol | Carisoprodol |
| Glyburide | Clonidine | Celebrex |
| Januvia | Digoxin | Clonazepam |
| Lantus | Diltiazem | Concerta |
| Lipitor | Diovan | Cyclobenzaprin |
| Lovastatin | Enalapril | Diazepam |
| Metformin | Flomax | Gabapentin |
| Niaspan | Furosemide | Hydrocodone |
| Pravastatin | Hydrochlorothiazide | Ibuprofen |
| Simvastatin | Isosorbide | Lorazepam |
| Tricor | Lisinopril | Lyrica |
| Vytorin | Metoprolol | Meloxicam |
| Zetia | Toprol | Namenda |
| | Triamterene | Naproxen |
| | Verapamil | Oxycodone |
| | | Oxycontin |
| | | Promethazine |
| | | Propoxphyene |
| | | Suboxone |
| | | Tramadol |

| | | Vyvanse |
| --- | --- | --- |
| | | Zolpidem |

**Table A.4 Listing of drugs that were classified as Metabolic Agents, Cardiovascular Agents, and Central Nervous System Agents.**

Figure A.1 visualizes the distribution of the drug categories listed in Tables A.1, A.2, and A.3. Central nervous system agents, cardiovascular agents, and metabolic agents attribute for roughly fifty percent of the drugs investigated in this work.



**Figure A.1 Distribution of drug categories for the list of drug names, as classified by the Drugs.com taxonomy.**

# Appendix B Statistics for each OSN

Table B.1 reports the number of posts, number of unique posts, posts per day, and the average length of each post. General OSNs tend to have many more duplicate posts than health OSNs due to advertisements and reposting of content. MediGuard is an exception since all drug reviews for a particular drug are listed under the brand name, and its search feature has up-to-date information on generic to brand drug name mappings.

General OSNs such as Twitter, Pinterest, and Google+ contain many more posts over a shorter period of time than the health OSNs. This difference is also emphasized by the number of posts per day. However, the average length of a post from a general OSN is much smaller than the average length of a post in all health OSNs, with the exception of DailyStrength. This is due to the nature of drug reviews in DailyStrength – a majority of reviews are short phrases such as "works for me" or "doesn't work".

| Dataset | Total Posts | Unique Posts | Percent Unique | Avg. Posts per Day | Std. Dev. | Avg. Words per Post | Std. Dev. |
|---|---|---|---|---|---|---|---|
| Twitter | 852,692 | 587,460 | 68.9% | 13,166* | 6,169* | 13.1 | 5.7 |
| Google+ | 11,803 | 8,706 | 73.7% | 15.5 | 25.7 | 39.6 | 70.1 |
| Pinterest | 8,706 | 5,876 | 66.5% | N/A | N/A | 24.9 | 33.2 |
| DailyStrength | 81,514 | 72,522 | 88.9% | N/A | N/A | 15.7 | 13.4 |
| Drugs.com | 5,451 | 4,994 | 91.6% | 2.4 | 2.3 | 64.5 | 41.8 |
| DrugLib.com | 974 | 959 | 98.4% | N/A | N/A | 121.4 | 84.1 |
| everydayHealth | 852 | 820 | 96.2% | 0.19 | 0.74 | 77.5 | 51.6 |
| MediGuard | 21,278 | 15,126 | 71.0% | 6.9 | 47.6 | 72.8 | 65.4 |
| medications | 35,050 | 34,997 | 99.8% | N/A | N/A | 133.9 | 135.6 |
| WebMD | 28,482 | 27,705 | 97.2% | 14.2 | 7.0 | 61.2 | 60.9 |

**Table B.1 General statistics for each of the OSNs. The total number of posts, total number of unique posts, average posts per day, and average words per post are given.**

**\* The average number of posts per day were computed based on all Tweets matching the drug name filter for the specified dates listed in Table A.1.**

Table B.2 summarizes the medical concept content of each OSN in terms of the number of medical concepts per post and per word. These results were computed across all concepts in a given post and for medical concepts that are unique in a given post. Except for DailyStrength, each health OSN contains a higher number of concepts per post, but the concentration of medical concepts per word is higher in general OSNs than health OSNs. These results coupled with the observations from Table B.1 suggest that users are sharing stories about their experiences with a particular drug in health OSNs; whereas users in general OSNs are expressing shorter thoughts

with more medical concepts, such as advertisements, news, educational material, or jokes. Again, the only

exception is DailyStrength.

| Dataset | Avg. Concepts per Post | Std. Dev. | Avg. Unique Concepts per Post | Std. Dev. | Avg. Concepts per Word | Avg. Unique Concepts per Word | Total Concepts | Unique Concepts |
|---|---|---|---|---|---|---|---|---|
| Twitter | 2.2 | 1.5 | 2.0 | 1.4 | 0.178 | 0.166 | 981,223 | 10,519 |
| Google+ | 6.1 | 8.4 | 5.3 | 5.9 | 0.181 | 0.165 | 53,218 | 5,849 |
| Pinterest | 4.8 | 5.8 | 4.3 | 4.6 | 0.219 | 0.201 | 28,136 | 4,067 |
| DailyStrength | 1.9 | 2.2 | 1.8 | 2.0 | 0.130 | 0.127 | 158,669 | 4,820 |
| Drugs.com | 9.7 | 6.6 | 8.3 | 5.3 | 0.158 | 0.142 | 48,425 | 3,500 |
| DrugLib.com | 19.6 | 12.3 | 15.0 | 7.8 | 0.173 | 0.138 | 18,818 | 2,305 |
| everydayHealth | 11.4 | 8.0 | 9.5 | 6.3 | 0.156 | 0.137 | 9,339 | 1,577 |
| MediGuard | 7.6 | 8.5 | 6.1 | 6.2 | 0.110 | 0.095 | 160,660 | 6,308 |
| Medications | 18.7 | 19.0 | 12.7 | 10.8 | 0.150 | 0.119 | 654,340 | 9,169 |
| WebMD | 8.8 | 8.6 | 7.5 | 6.5 | 0.167 | 0.153 | 244,589 | 6,535 |

**Table B.2 Overview of medical concept content. The average number of concepts, total number of concepts,**

**and the average number of concepts per word are shown; these results only consider concepts from semantic**

**groups related to medicine.**

Tables B.3 and B.4 summarize the distribution of the drug categories for each OSN. Each post in a given

OSN was assigned to a drug from Tables A.2-A.4. The number of posts for each drug category was then tallied

and the distribution of a given drug category was calculated by dividing each tally by the total number of posts for

the given OSN. These tables show trends amongst OSNs, such as the dominance of genitourinary tract agents and

nutritional products amongst general OSNs, and the dominance of psychotherapeutic agents amongst health OSNs.

| | Alternative Medicines | Anti-infectives | Cardiovascular Agents | Central Nervous System Agents | Coagulation Modifiers | Gastrointestinal Agents |
|---|---|---|---|---|---|---|
| Twitter | 0.10% | 6.21% | 5.05% | 34.69% | 0.98% | 1.50% |
| Google+ | 0.13% | 6.10% | 4.40% | 24.71% | 3.91% | 1.27% |
| Pinterest | 0.09% | 12.63% | 2.74% | 28.18% | 1.55% | 1.34% |
| DailyStrength | 0.00% | 0.77% | 2.57% | 23.72% | 0.14% | 0.59% |
| Drugs.com | 0.08% | 8.75% | 11.13% | 28.38% | 0.96% | 2.42% |
| DrugLib.com | 0.00% | 5.63% | 3.86% | 25.55% | 0.52% | 4.38% |
| everydayHealth | 0.00% | 5.37% | 11.46% | 32.01% | 0.85% | 1.22% |
| MediGuard | 0.00% | 1.39% | 14.48% | 23.02% | 2.19% | 4.14% |
| Medications | 0.00% | 14.78% | 19.72% | 2.30% | 0.13% | 1.77% |
| WebMD | 0.69% | 2.73% | 20.31% | 28.33% | 1.23% | 1.37% |

**Table B.3 Distribution of the drug categories for each OSN.**

| | Genitourinary Track Agents | Hormones | Metabolic Agents | Nutritional Products | Psychotherapeutic Agents | Respiratory Agents | Topical Agents |
|---|---|---|---|---|---|---|---|
| Twitter | 29.03% | 2.87% | 6.87% | 1.27% | 5.51% | 5.13% | 0.78% |
| Google+ | 13.65% | 6.51% | 16.88% | 6.85% | 8.72% | 6.42% | 0.47% |
| Pinterest | 8.29% | 2.60% | 14.25% | 15.89% | 3.81% | 8.07% | 0.56% |
| DailyStrength | 0.38% | 8.52% | 6.26% | 0.03% | 50.04% | 6.64% | 0.32% |
| Drugs.com | 1.30% | 7.21% | 5.81% | 0.16% | 28.21% | 5.02% | 0.56% |
| DrugLib.com | 1.25% | 14.49% | 10.79% | 0.57% | 26.49% | 3.96% | 2.50% |
| everydayHealth | 1.34% | 2.99% | 8.29% | 0.00% | 33.78% | 2.26% | 0.43% |
| MediGuard | 0.85% | 10.81% | 15.82% | 1.69% | 18.68% | 6.13% | 0.79% |
| Medications | 0.10% | 25.71% | 14.59% | 0.03% | 3.09% | 17.36% | 0.42% |
| WebMD | 1.15% | 6.55% | 12.43% | 0.05% | 23.66% | 1.10% | 0.40% |

**Table B.4 Distribution of the drug categories for each OSN.**

Table B.5 reports the average sentiment for each OSN. This table shows the trend that health OSNs tend to be more subjective than general OSNs, with the exception of MediGuard.

| | Positive | Negative | Objective |
|---|---|---|---|
| Twitter | 5.73% | 6.50% | 87.77% |
| Google+ | 2.91% | 2.66% | 94.43% |
| Pinterest | 6.31% | 5.51% | 88.18% |
| DailyStrength | 7.33% | 9.46% | 83.21% |
| Drugs.com | 7.51% | 8.44% | 84.05% |
| DrugLib.com | 6.55% | 8.37% | 85.07% |
| everydayHealth | 6.59% | 8.56% | 84.85% |
| MediGuard | 3.20% | 3.52% | 93.29% |
| Medications | 5.86% | 8.01% | 86.12% |
| WebMD | 6.66% | 8.67% | 84.67% |

**Table B.5 Distribution of the average sentiment for each OSN.**

Table B.6 reports the distribution of semantic groups for each OSN. This table shows that disorders and physiology are mentioned more often in health OSNs, whereas general OSNs mentioned explicit drug names.

| | Anatomy | Chemicals and Drugs | Physiology | Disorders | Procedures |
|---|---|---|---|---|---|
| Twitter | 4.98% | 44.13% | 12.64% | 25.55% | 12.70% |
| Google+ | 6.28% | 37.15% | 11.59% | 29.85% | 15.13% |
| Pinterest | 6.37% | 45.57% | 10.08% | 24.79% | 13.18% |

| | | | | | |
|---|---|---|---|---|---|
| DailyStrength | 5.07% | 13.76% | 21.78% | 47.50% | 11.90% |
| Drugs.com | 6.49% | 19.49% | 17.22% | 45.54% | 11.26% |
| DrugLib.com | 6.89% | 17.17% | 18.66% | 42.29% | 14.99% |
| everydayHealth | 6.68% | 19.05% | 17.26% | 43.62% | 13.38% |
| MediGuard | 7.47% | 23.24% | 16.20% | 37.34% | 15.75% |
| Medications | 11.14% | 15.15% | 19.03% | 42.96% | 11.72% |
| WebMD | 7.82% | 17.27% | 18.95% | 43.58% | 12.38% |

**Table B.6 Distribution of semantic groups for each OSN.**

Table B.7 reports the distribution of emotion for each OSN. Overall the results are not as interesting as the pervious tables, with the exception of the number of trust terms contained in MediGuard.

| | Anger | Fear | Trust | Disgust | Anticipation | Surprise |
|---|---|---|---|---|---|---|
| Twitter | 39.45% | 60.55% | 64.28% | 35.72% | 75.28% | 24.72% |
| Google+ | 33.34% | 66.66% | 68.92% | 31.08% | 77.39% | 22.61% |
| Pinterest | 31.80% | 68.20% | 69.15% | 30.85% | 73.54% | 26.46% |
| DailyStrength | 30.17% | 69.83% | 57.07% | 42.93% | 73.61% | 26.39% |
| Drugs.com | 30.33% | 69.67% | 68.66% | 31.34% | 70.41% | 29.59% |
| DrugLib.com | 34.29% | 65.71% | 65.36% | 34.64% | 75.05% | 24.95% |
| everydayHealth | 32.52% | 67.48% | 66.49% | 33.51% | 71.64% | 28.36% |
| MediGuard | 28.09% | 71.91% | 74.31% | 25.69% | 72.89% | 27.11% |
| Medications | 34.94% | 65.06% | 60.74% | 39.26% | 72.52% | 27.48% |
| WebMD | 30.55% | 69.45% | 64.50% | 35.50% | 72.26% | 27.74% |

**Table B.7 Distribution of emotional pairs for each OSN.**

# Appendix C General versus Health OSNs

Table C.1 summarizes general and medical concept statistics for the two groupings of OSNs. General OSNs contain more posts with fewer words per post, but general OSNs have a smaller percentage of unique posts. Health OSNs contain more concepts per post due to their increased length, but these OSNs have fewer concepts per word.

| Category | Total Posts | Unique Posts | Words Per Post | Average Concepts per Post | Avg. Concepts per Word | Unique Concepts |
|---|---|---|---|---|---|---|
| General | 873,201 | 602,042 (69%) | 25.9 | 4.4 | 0.177 | 13,238 |
| Health | 173,601 | 157,123 (90%) | 64.7 | 9.7 | 0.149 | 13,130 |

**Table C.1 Summary of general statistics and medical concept statistics for general and health OSNs.**

Figure C.1 compares distributions of emotional pairs of the health and general OSNs with the distribution of a uniform baseline, where the baseline assumes a uniform distribution for every term mapped from the NRC word-

emotion lexicon [3]. This lexicon contains over 14,000 words manually labeled by humans via crowdsourcing. Each term is assigned one or more emotional-pairs from the following set: (1) negative–positive; (2) joy–sadness; (3) anger–fear; (4) trust–disgust; and (5) anticipation–surprise. Since joy–sadness is similar to positive–negative, and we compute positive, negative, and objective scores using SentiWordNet [4], our analysis omits results for the emotional pairs joy–sadness and positive–negative from the NRC lexicon. Analogous to the SentiWordNet process, we stemmed both the posts and the terms in the NRC lexicon before computing the emotion scores. We then mapped phrases from the NRC lexicon to phrases in the posts using the longest possible match first. Next, we computed the score for each emotional-pair of each post by averaging the emotion scores from every mapped term. The final score for each emotional pair is then computed by averaging the emotion scores of all posts within a given OSN. Table C.2 reports the highest absolute relative changes of each emotional pair shown in Figure C.1 relative to the baseline. All comparisons in this figure are with $p < 0.001$ both significance tests. Health and general OSNs follow the same trends with respect to the baseline; both groups observe an increase in fear, trust, and anticipation terms, and a decrease in anger, disgust, and surprise terms.



**Figure C.1 An overview of the emotion analysis for general OSNs versus health OSNs. (A) The distribution of fear–anger; (B) the distribution of disgust–trust; and (C) the distribution of surprise–anticipation.**

| Emotion | | |
|---|---|---|
| Surprise | General | -36% |
| Anger | Health | -31% |
| Surprise | Health | -29% |
| Disgust | General | -29% |
| Fear | Health | +26% |
| Trust | General | +25% |
| Disgust | Health | -24% |
| Anger | General | -23% |
| Anticipation | General | +23% |

| Trust | Health | +21% |
| Fear | General | +20% |
| Anticipation | Health | +18% |

**Table C.2 Highest absolute relative changes of the emotional pairs compared with the baselines shown in Figure C.1. E.g., General Surprise is computed as the difference between General Surprise and Baseline Surprise divided by Baseline Surprise.**

Frequent itemsets of medical concepts for health and general OSNs are reported in Tables C.3 and C.4. General OSNs are dominated by specific drug names, where drugs with similar purposes often co-occur together, such as Ibuprofen, Tylenol, and Advil. Drugs co-occur in a single post for multiple reasons:

- Online pharmacies advertise multiple drugs that serve a single purpose; e.g., "[URL] order Viagra Cialis and Levitra in internet shop without script California !".

- Users will associate conditions with each drug from a single group; e.g., "WORSE HEADACHE EVER!!! #TYLENOL #IBUPROFEN #ADVIL"

Drugs such as Viagra are also often used in jokes; e.g. and "Viagra for women has been around for centuries. It's called money".

Another interesting itemset from Table C.4 is Viagra, watching, and awkward; this itemset refers to posts that discuss the awkwardness of watching Viagra commercials with one's family. Lastly, there are a series of itemsets referring to Viagra, death, overdose, and amputated. These itemsets are referring to jokes or odd news articles that refer to the comical effects of taking too much Viagra.

As shown in Table C.3, sleep occurs in six of the ten most frequent itemsets of size two. Several itemsets from Tables C.3 and C.4 refer to drugs and the conditions they treat:

- Lisinopril and hypertension.

- Singular and asthma.

- Sleep and Ambien.

- Lexapro, anxiety, and depression

Itemsets of symptoms are also common to health OSNs, such as headache, dizziness, and nausea. Health OSNs also contain frequent itemsets of drugs and their side effects: Lisinopril, sleepiness, and coughing; NuvaRing and decreased Libido; Singulair and depression.

We further examined frequent itemsets of all possible concepts for both general and health OSNs, reported in Tables C.5 and C.6. These itemsets yield further insight into the types of conversations users have in each grouping of OSNs. Several itemsets for general OSNs are related to advertisements from online pharmacies; these itemsets include concepts such as Internet, mail, priority, prices, low, scripts, and order.

Itemsets for health OSNs reveal that users are discussing their experiences with their medications, and the differing strategies employed by their physicians; the concept for physicians appears in over half of the frequent itemsets for both Tables C.5 and C.6. These posts typically discuss a problem and an action; e.g., "my doctor increased my dosage to 20mg"; and "[My doctor] put me on Lisinopril but stopped taking it after constantly coughing day and night".

| Health OSNs | | | General OSNs | | |
|---|---|---|---|---|---|
| Depression | Anxiety | 1.05% | Viagra | Sexual Intercourse | 2.24% |
| Lisinopril | Coughing | 0.91% | Ibuprofen | Headache | 1.06% |
| Singulair | Asthma | 0.79% | Online Pharmaceutical Services | Buying drugs | 0.72% |
| Sleep | Tired | 0.75% | Ibuprofen | Acetaminophen | 0.65% |
| Sleep | Sleeplessness | 0.65% | Viagra | Male | 0.60% |
| Sleep | Depression | 0.64% | Ibuprofen | Ice | 0.55% |
| Sleep | Eating | 0.61% | Viagra | Penile Erection | 0.46% |
| Sleep | Anxiety | 0.60% | Ibuprofen | Sleep | 0.39% |
| Headache | Nausea | 0.58% | Viagra | Female | 0.34% |
| Sleep | Ambien | 0.57% | Ibuprofen | Tylenol | 0.33% |

**Table C.3 Frequent itemsets of size 2 for medical concepts.**

| Health OSNs | | | | General OSNs | | | |
|---|---|---|---|---|---|---|---|
| Singulair | Hypersensitivity | Asthma | 0.25% | Acids | Abdominal Colic | Autistic | 0.11% |
| Sleep | Lisinopril | Coughing | 0.16% | Viagra | Decision | Female | 0.07% |
| Lisinopril | Blood pressure | Coughing | 0.16% | Acids | Abdominal Colic | Pregnancy | 0.07% |
| Sleep | Depression | Anxiety | 0.15% | Viagra | Watching | Awkward | 0.06% |
| Depression | Anxiety | Lexapro | 0.14% | Viagra | Overdose | Death | 0.05% |
| Libido | NuvaRing | Sexual Intercourse | 0.14% | Ibuprofen | Tylenol | Advil | 0.05% |
| Headache | Dizziness | Nausea | 0.13% | Amphetamine | Withdrawal | Amphetamine Withdrawal | 0.06% |
| Depression | Singulair | Asthma | 0.12% | Viagra | Penis | Amputated | 0.05% |
| Libido | NuvaRing | Contraceptives | 0.11% | Coughing | Cough Syrup | Codeine | 0.05% |
| Lisinopril | Blood pressure | Hypertension | 0.11% | Viagra | Overdose | Amputated | 0.05% |

**Table C.4 Frequent itemsets of size 3 for medical concepts.**

| Health OSNs | | | General OSNs | | |
|---|---|---|---|---|---|
| Physicians | Started | 3.83% | Viagra | Sexual Intercourse | 1.77% |
| Physicians | Milligram | 3.73% | Viagra | Prices | 1.57% |
| Milligram | Started | 3.34% | Prices | Lowest | 1.42% |
| Help | Physicians | 3.02% | Viagra | Lowest | 1.39% |
| Milligram | Dosage | 3.01% | Overnight | Transfer | 1.12% |
| Help | Sleep | 2.69% | Viagra | Commercial | 0.96% |
| Help | Milligram | 2.63% | Ibuprofen | Headache | 0.90% |
| Physicians | Dosage | 2.55% | Order | Scripts | 0.78% |
| Physicians | Better | 2.25% | Order | Internet | 0.72% |
| Help | Started | 2.14% | Ibuprofen | Milligram | 0.60% |

**Table C.5 Frequent itemsets of size 2 for all UMLS concepts.**

| Health OSNs | | | | General OSNs | | | |
|---|---|---|---|---|---|---|---|
| Physicians | Milligram | Started | 1.31% | Viagra | Prices | Lowest | 1.38% |
| Physicians | Milligram | Dosage | 1.17% | Order | Overnight | Transfer | 0.44% |
| Milligram | Started | Dosage | 0.92% | Scripts | Overnight | Transfer | 0.14% |
| Help | Physicians | Milligram | 0.89% | Viagra | Commercial | Watching | 0.14% |
| Help | Physicians | Started | 0.88% | Viagra | Commercial | Awkward | 0.11% |
| Physicians | Started | Better | 0.84% | Pharmacy | Overnight | Transfer | 0.11% |
| Physicians | Started | Dosage | 0.78% | Order | Overnight | Delivery | 0.11% |
| Physicians | Milligram | Better | 0.72% | Viagra | Overnight | Delivery | 0.11% |
| Help | Milligram | Started | 0.70% | Viagra | Commercial | Awkward | 0.10% |
| Physicians | Started | Last | 0.69% | Order | Internet | Scripts | 0.10% |

**Table C.6 Frequent itemsets of size 3 for all UMLS concepts.**

# Appendix D Non-moderated versus Moderated Health OSNs

Table D.1 summarizes general and medical concept statistics for moderated and non-moderated health OSNs. Moderated OSNs contain many more words per post, due to their inclusion of medications and DrugLib.com, both of which contain over 120 words per post. Thus, moderated health OSNs also contain more concepts per post and cover more concepts than non-moderated health OSNs.

| Category | Total Posts | Unique Posts | Words Per Post | Average Concepts per Post | Avg. Concepts per Word | Unique Concepts |
|---|---|---|---|---|---|---|
| Not Moderated | 110,848 | 101,047 (91%) | 51.5 | 7.4 | 0.151 | 7,875 |
| Moderated | 62,753 | 56,076 (89%) | 99.3 | 13.9 | 0.148 | 11,651 |

**Table D.1 Summary of general statistics and medical concept statistics for not moderated and moderated OSNs.**

Figure D.1 reports the effect of moderation on the emotional pairs; all comparisons in this figure are significant with p < 0.001 for both significance tests, except for the Pearson's Chi Squared test of independence for anticipation

and surprise. Moderated OSNs decreased the number of disgusting terms and increased the number of trusting terms, whereas lack of moderation had the opposite effect. Otherwise, moderation had little or no effect on the emotional and medical content of drug reviews in health OSNs.



**Figure D.1 An overview of the emotion analysis for moderated and not moderated OSNs. (A)-(C) the distribution of the emotional pairs fear–anger, disgust–trust, and surprise–anticipation.**

| Emotion | | |
|---|---|---|
| Disgust | Not Moderated | +13% |

**Table D.2 Highest absolute relative changes of an emotion for a given OSN group compared with the emotion of the other OSN grouping. E.g., Not Moderated Disgust is computed as the difference between Not Moderated Disgust and Moderated Disgust divided by Moderated Disgust.**

Tables D.3-5 report frequent itemsets for health OSNs with and without moderation. Sleep is common to both groupings of OSNs, but sleep is more frequent for non-moderated health OSNs. Frequent itemsets from non-moderated health OSNs concur with Figure 4(A) from Section 4.2, in that psychotherapeutic agents (Lexapro and Cymbalta), along with psychological conditions (panic attacks, mental suffering, depression, and anxiety) are frequent; whereas these drugs are not found in the frequent itemsets of moderated health OSNs, and these conditions are not as frequent. Moderated health OSNs contain concepts relating to the respiratory and cardiovascular systems, including Lisinopril, Singulair, Lipitor, asthma, coughing, hypertension, blood pressure, and cholesterol. Further, moderated health OSNs also discuss the contraceptive NuvaRing and its side effect of decreased libido.

| Non-moderated Health OSNs | | Moderated Health OSNs | |
|---|---|---|---|
| Sleep | 10.93% | Sleep | 9.19% |
| Depression | 5.11% | Lisinopril | 6.44% |

| | | | | |
|---|---|---|---|
| Weight Gain | 3.70% | Singulair | 6.25% |
| Tired | 3.70% | Depression | 6.01% |
| Anxiety | 3.63% | Headache | 5.66% |
| Headache | 3.00% | Mental Suffering | 5.61% |
| Dizziness | 2.68% | Eating | 5.26% |
| Drowsiness | 2.58% | Prednisone | 4.91% |
| Nausea | 2.52% | Levaquin | 4.61% |
| Lexapro | 2.33% | Tired | 4.47% |

**Table D.3 Frequent itemsets of size 1 for medical concepts.**

| Non-moderated Health OSNs | | | Moderated Health OSNs | | |
|---|---|---|---|---|---|
| Depression | Anxiety | 1.10% | Lisinopril | Coughing | 1.82% |
| Sleep | Anxiety | 0.57% | Singulair | Asthma | 1.81% |
| Sleep | Depression | 0.57% | Lisinopril | Listerine | 1.51% |
| Sleep | Ambien | 0.56% | Singulair | Hypersensitivity | 1.25% |
| Depression | Lexapro | 0.55% | Lipitor | Cholesterol | 1.12% |
| Sleep | Tired | 0.52% | Sleep | Tired | 1.06% |
| Depression | Cymbalta | 0.51% | Lisinopril | Blood pressure | 1.04% |
| Sleep | Sleeplessness | 0.50% | Sleep | Depression | 1.00% |
| Sleep | Eating | 0.46% | Depression | Anxiety | 0.98% |
| Sleep | Drowsiness | 0.41% | Asthma | Hypersensitivity | 0.91% |

**Table D.4 Frequent itemsets of size 2 for medical concepts.**

| Non-moderated Health OSNs | | | | Moderated Health OSNs | | | |
|---|---|---|---|---|---|---|---|
| Depression | Anxiety | Lexapro | 0.19% | Singulair | Asthma | Hypersensitivity | 0.58% |
| Sleep | Depression | Anxiety | 0.14% | Sleep | Lisinopril | Coughing | 0.34% |
| Depression | Anxiety | Cymbalta | 0.12% | NuvaRing | Libido | Sexual Intercourse | 0.32% |
| Headache | Dizziness | Nausea | 0.09% | Lisinopril | Coughing | Blood pressure | 0.31% |
| Sleep | Depression | Sleeplessness | 0.08% | Singulair | Depression | Asthma | 0.27% |
| Sleep | Sleeplessness | Ambien | 0.08% | Singulair | Happiness | Asthma | 0.25% |
| Depression | Anxiety | Panic Attacks | 0.08% | Singulair | Mental Suffering | Asthma | 0.24% |
| Depression | Anxiety | Mental Suffering | 0.07% | NuvaRing | Libido | Contraceptives | 0.23% |
| Depression | Weight Gain | Anxiety | 0.07% | Sleep | Singulair | Asthma | 0.23% |
| Sleep | Remembering | Ambien | 0.07% | Lisinopril | Blood pressure | Hypertension | 0.23% |

**Table D.5 Frequent itemsets of size 3 for medical concepts.**

## Appendix E Registration versus No Registration for Health OSNs

Table E.1 summarizes general and medical concept statistics for health OSNs that do or do not require registration. Registration has little effect on these statistics, with the average number of words and medical concepts being roughly equal.

| Category | Total Posts | Unique Posts | Words Per Post | Average Concepts per Post | Avg. Concepts per Word | Unique Concepts |
|---|---|---|---|---|---|---|
| No Registration | 35,759 | 34,478 (96%) | 81.2 | 12.4 | 0.164 | 7,567 |

| Registration | 137,842 | 122,645 (89%) | 74.3 | 9.4 | 0.130 | 11,839 |

**Table E.1 Summary of general statistics and medical concept statistics for health OSNs that do and do not require registration.**

Figure E.1 reports the effect of registration on the emotional pairs; all comparisons in this figure are significant with p < 0.001 for both significance tests. Overall, registration had little or no effect on the emotional and medical content of drug reviews in health OSNs.



**Figure E.1 An overview of the emotional analysis for health OSNs that do and do not require registration.**

**(A)-(C) The distribution of the emotional pairs fear–anger, disgust–trust, and surprise–anticipation.**

| Emotion | | |
|---|---|---|
| Disgust | Registration | +6% |

**Table E.2 Highest absolute relative changes of an emotion for a given OSN group compared with the emotion of the other OSN grouping. E.g., Registration Disgust is computed as the difference Registration Disgust and No Registration Disgust divided by No Registration Disgust.**

Tables D.3-5 report frequent itemsets for health OSNs that do or do not require registration. Similar to moderation, sleep is common to both groupings of OSNs, but is more prevalent in health OSNs that do not require registration. Further, concepts relating to psychotherapeutics and psychological conditions are common in health OSNs that do not require registration; analogous to the frequent itemsets for non-moderated health OSNs. Health OSNs that require registration have similar frequent itemsets to that of health OSNs with moderation, which focus on respiratory and cardiovascular drugs and conditions, such as Lisinopril, Lipitor, Singulair, asthma, and allergies. Also similar to health OSNs with moderation, health OSNs with registration have NuvaRing and its side effect libido as frequent itemsets.

| No Registration Health OSNs | | Registration Health OSNs | |
| --- | --- | --- | --- |
| Sleep | 12.26% | Sleep | 9.51% |
| Depression | 6.51% | Depression | 4.61% |
| Headache | 5.59% | Vision | 3.93% |
| Tired | 5.36% | Weight Gain | 3.76% |
| Anxiety | 4.92% | Headache | 3.62% |
| Dizziness | 4.83% | Personal appearance | 3.61% |
| Happiness | 4.49% | Lisinopril | 3.58% |
| Nausea | 4.42% | Tired | 3.58% |
| Blood pressure finding | 4.28% | Singulair | 3.48% |
| Weight Gain | 4.18% | Eating | 3.29% |

**Table E.3 Frequent itemsets of size 1 for medical concepts.**

| No Registration Health OSNs | | | Registration Health OSNs | | |
| --- | --- | --- | --- | --- | --- |
| Depression | Anxiety | 1.78% | Lisinopril | Coughing | 1.01% |
| Sleep | Ambien | 1.19% | Singulair | Asthma | 0.99% |
| Sleep | Depression | 1.16% | Depression | Anxiety | 0.80% |
| Depression | Lexapro | 1.13% | Singulair | Hypersensitivity | 0.69% |
| Sleep | Sleeplessness | 1.13% | Sleep | Tired | 0.65% |
| Sleep | Tired | 1.05% | Lipitor | Cholesterol | 0.60% |
| Sleep | Anxiety | 1.01% | Lisinopril | Blood Pressure | 0.57% |
| Depression | Cymbalta | 0.99% | Sleep | Depression | 0.57% |
| Sleep | Eating | 0.94% | Hypersensitivity | Asthma | 0.53% |
| Dizziness | Nausea | 0.85% | Headache | Nausea | 0.50% |

**Table E.4 Frequent itemsets of size 2 for medical concepts.**

| No Registration Health OSNs | | | | Registration Health OSNs | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| Depression | Anxiety | Lexapro | 0.41% | Singulair | Hypersensitivity | Asthma | 0.32% |
| Sleep | Depression | Anxiety | 0.38% | Sleep | Lisinopril | Coughing | 0.19% |
| Depression | Anxiety | Cymbalta | 0.25% | NuvaRing | Libido | Sexual Intercourse | 0.18% |
| Sleep | Depression | Sleeplessness | 0.24% | Lisinopril | Coughing | Blood pressure | 0.17% |
| Depression | Anxiety | Happiness | 0.23% | Depression | Singulair | Asthma | 0.15% |
| Headache | Dizziness | Nausea | 0.22% | Singulair | Happiness | Asthma | 0.14% |
| Depression | Anxiety | Panic Attacks | 0.21% | Singulair | Suffering | Asthma | 0.13% |
| Sleep | Sleeplessness | Ambien | 0.20% | NuvaRing | Libido | Contraceptives | 0.13% |
| Depression | Anxiety | Weight Gain | 0.19% | Sleep | Singulair | Asthma | 0.12% |
| Depression | Anxiety | Zoloft | 0.19% | NuvaRing | Sexual Intercourse | Contraceptives | 0.12% |

**Table E.5 Frequent itemsets of size 3 for medical concepts.**

# Appendix F Review versus Q&A OSNs

Table F.1 summarizes general and medical concept statistics for health OSNs with a review format and with a Q&A format. The format has little effect on these statistics, with the average number of words and medical concepts being roughly equal.

| Category | Total Posts | Unique Posts | Words Per Post | Average Concepts per Post | Avg. Concepts per Word | Unique Concepts |
| --- | --- | --- | --- | --- | --- | --- |
| Review | 152,323 | 141,997 (93%) | 62.4 | 11.6 | 0.156 | 12,152 |

| Q&A | 21,278 | 15,126 (71%) | 72.8 | 7.6 | 0.110 | 6,308 |

**Table F.1 Summary of general statistics and medical concept statistics for health OSNs with a review format and Q&A format.**

Figure F.1 reports the effect of health OSN format on the emotional pairs; all comparisons in F.1(A) and F.1(B) are significant with p < 0.001 for both significance tests. Health OSNs with a Q&A format observed a 26% decrease in disgusting terms and a 20% increase in trusting terms. Further, the format has little effect on the emotional pair surprise–anticipation; however, health OSNs with a Q&A format observed a decrease of 10% in anger terms, and an increase of 5% in fear terms.



**Figure F.1 An overview of the emotional analysis for health OSNs with a review format and a Q&A format. (A) The distribution of fear–anger,; (B) the distribution of surprise–anticipation.**

| Emotion | | |
|---|---|---|
| Disgust | Review | +40% |
| Trust | Q&A | +16% |
| Anger | Review | +14% |

**Table F.2 Highest absolute relative changes of an emotion for a given OSN group compared with the emotion of the other OSN grouping. E.g., Review Disgust is computed as the difference Review Disgust and Q&A Disgust divided by Q&A Disgust**

Sleep, anxiety, and depression are common and prevalent amongst both groupings of OSNs. Lisinopril, Lipitor, and NuvaRing are observed in health OSNs with a review format, but not those with a Q&A format. Lastly, there are frequent itemsets related to Xanax, Zoloft, hypothyroidism, and Synthroid in health OSNs with a Q&A format.

| Review Format Health OSNs | | Q&A Format Health OSNs | |
|---|---|---|---|
| Sleep | 10.47% | Female | 9.02% |
| Depression | 5.07% | Sleep | 8.04% |

| Headache | 4.32% | Eating | 4.95% |
|---|---|---|---|
| Tired | 4.17% | Disease | 4.36% |
| Weight Gain | 3.82% | Weight Gain | 4.26% |
| Anxiety | 3.67% | Anxiety | 4.24% |
| Lisinopril | 3.56% | Male gender | 4.04% |
| Eating | 3.45% | Depression | 3.60% |
| Dizziness | 3.35% | Mental Suffering | 3.53% |
| Nausea | 3.29% | Thyroid Gland | 3.50% |

**Table F.3 Frequent itemsets of size 1 for medical concepts.**

| Review Format Health OSNs | | | Q&A Format Health OSNs | | |
|---|---|---|---|---|---|
| Depression | Anxiety | 1.04% | Sleep | Ambien | 1.18% |
| Lisinopril | Coughing | 1.01% | Anxiety | Depression | 1.11% |
| Singulair | Asthma | 0.85% | Thyroid Gland | Synthroid | 1.10% |
| Sleep | Tired | 0.76% | Female | Sleep | 0.89% |
| Sleep | Sleeplessness | 0.71% | Sleep | Anxiety | 0.80% |
| Headache | Nausea | 0.65% | Anxiety | Panic Attacks | 0.72% |
| Sleep | Eating | 0.63% | Anxiety | Xanax | 0.71% |
| Sleep | Depression | 0.62% | Sleep | Tired | 0.70% |
| Singulair | Hypersensitivity | 0.59% | Sleep | Xanax | 0.70% |
| Lipitor | Cholesterol | 0.59% | Sleep | Sleeplessness | 0.67% |

**Table F.4 Frequent itemsets of size 2 for medical concepts.**

| Review Format Health OSNs | | | | Q&A Format Health OSNs | | | |
|---|---|---|---|---|---|---|---|
| Singulair | Hypersensitivity | Asthma | 0.27% | Thyroid Gland | Synthroid | Hypothyroidism | 0.30% |
| Sleep | Lisinopril | Coughing | 0.18% | Anxiety | Depression | Zoloft | 0.23% |
| Lisinopril | Coughing | Blood Pressure | 0.17% | Sleep | Anxiety | Depression | 0.19% |
| NuvaRing | Libido | Sexual Intercourse | 0.16% | Sleep | Anxiety | Xanax | 0.19% |
| Sleep | Depression | Anxiety | 0.15% | Sleep | Sleeplessness | Ambien | 0.19% |
| Depression | Anxiety | Lexapro | 0.15% | Disease | Thyroid Gland | Synthroid | 0.18% |
| Headache | Dizziness | Nausea | 0.15% | Weight Gain | Thyroid Gland | Synthroid | 0.18% |
| Depression | Singulair | Asthma | 0.13% | Anxiety | Depression | Celexa | 0.18% |
| NuvaRing | Libido | Contraceptives | 0.12% | Thyroid Gland | Synthroid | Blood | 0.17% |
| Lisinopril | Coughing | Dry cough | 0.12% | Female | Thyroid Gland | Synthroid | 0.16% |

**Table F.5 Frequent itemsets of size 3 for medical concepts.**

## Appendix G Statistical Tests

Tables G.1 and G.2 report the p-values for Pearson's Chi Squared test of independence and the Mann-Whitney U test. Note that all interpretations are based on statistically significant results.

| | Drug Categories | Sentiment | Semantic Groups | Anger – Fear | Trust – Disgust | Anticipation – Surprise |
|---|---|---|---|---|---|---|
| General vs Baseline | < 0.001 Figure 3(A) | < 0.001 Figure 3(B) | < 0.001 Figure 3(C) | < 0.0001 Figure C.1(A) | < 0.0001 Figure C.1(B) | < 0.001 Figure C.1(C) |

| Health vs Baseline | < 0.001 Figure 3(A) | < 0.001 Figure 3(B) | < 0.001 Figure 3(C) | < 0.001 Figure C.1(A) | < 0.001 Figure C.1(B) | < 0.001 Figure C.1(C) |
|---|---|---|---|---|---|---|
| General vs Health | < 0.001 Figure 3(A) | < 0.001 Figure 3(B) | < 0.001 Figure 3(C) | < 0.001 Figure C.1(A) | < 0.001 Figure C.1(B) | < 0.001 Figure C.1(C) |
| Moderated vs Not Moderated | < 0.001 Figure 4(A) | < 0.001 Figure 4(B) | < 0.001 Figure 4(C) | < 0.001 Figure D.1(A) | < 0.001 Figure D.1(B) | 0.013 Figure D.1(C) |
| Registration vs No Registration | < 0.001 Figure 5(A) | < 0.001 Figure 5(B) | < 0.001 Figure 5(C) | < 0.001 Figure E.1(A) | < 0.001 Figure E.1(B) | < 0.001 Figure E.1(C) |
| Review vs Q&A | < 0.001 Figure 6(A) | < 0.001 Figure 6(B) | < 0.001 Figure 6(C) | < 0.001 Figure F.1(A) | < 0.001 Figure F.1(B) | 0.01 Figure F.1(C) |

**Table G.1 p-values for Pearson's Chi Squared test of independence.**

| | General versus Health | Moderated vs Not Moderated | Registration vs No Registration | Review vs Q&A |
|---|---|---|---|---|
| Alternative Medicines | 0.04 | < 0.001 | < 0.001 | < 0.001 |
| Anti-infectives | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Cardiovascular Agents | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Central Nervous System Agents | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Coagulation Modifiers | < 0.001 | < 0.001 | 0.4345 | < 0.001 |
| Gastrointestinal Agents | 0.02 | < 0.001 | < 0.001 | < 0.001 |
| Genitourinary Track Agents | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Hormones | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Metabolic Agents | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Nutritional Products | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Psychotherapeutic Agents | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Respiratory Agents | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Topical Agents | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Positive | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Negative | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Objective | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Disorders | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Procedures | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Anatomy | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Drugs | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Physiology | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Anger | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Fear | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Trust | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Disgust | < 0.001 | < 0.001 | < 0.001 | < 0.001 |
| Anticipation | < 0.001 | < 0.001 | < 0.001 | 0.5874 |

| | | | |
|---|---|---|---|
| Surprise | < 0.001 | < 0.001 | < 0.001 | 0.5874 |

**Table G.2 p-values for Mann-Whitney U test on each post for a given OSN grouping. The test is computed for each variable (e.g., Alternative Medicines) from each category (e.g., drug categories).**

**References**

[1] RxList. RxList - The Interent Drug Index. http://www.rxlist.com/script/main/hp.asp (Accessed: January 15, 2013).

[2] Drugs.com. Drugs By Category. http://www.drugs.com/drug-classes.html?tree=1 (Accessed: March 4, 2013).

[3] Mohammad SM, Turney PD. Crowdsourcing a word–emotion association lexicon. Comput Intell. 2012. http://dx.doi.org/10.1111/j.1467-8640.2012.00460.x

[4] Baccianella S, Esuli A, Sebastiani F. Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10). Valletta, Malta: European Language Resources Association (ELRA); 2010.