<u>Original Paper</u>

# Deep Learning for Identification of Alcohol-Related Content on Social Media (Reddit and Twitter): Exploratory Analysis of Alcohol-Related Outcomes

Benjamin Joseph Ricard[1], MSc; Saeed Hassanpour[1,2,3], PhD

[1]Department of Biomedical Data Science, Dartmouth College, Lebanon, NH, United States

[2]Department of Epidemiology, Dartmouth College, Hanover, NH, United States

[3]Department of Computer Science, Dartmouth College, Hanover, NH, United States

**Corresponding Author:**
Saeed Hassanpour, PhD
Department of Biomedical Data Science
Dartmouth College
Williamson Translational Research Building
One Medical Center Drive HB 7261
Lebanon, NH, 03756
United States
Phone: 1 603 650 1983
Email: saeed.hassanpour@dartmouth.edu

## *Abstract*

**Background:** Many social media studies have explored the ability of thematic structures, such as hashtags and subreddits, to identify information related to a wide variety of mental health disorders. However, studies and models trained on specific themed communities are often difficult to apply to different social media platforms and related outcomes. A deep learning framework using thematic structures from Reddit and Twitter can have distinct advantages for studying alcohol abuse, particularly among the youth in the United States.

**Objective:** This study proposes a new deep learning pipeline that uses thematic structures to identify alcohol-related content across different platforms. We apply our method on Twitter to determine the association of the prevalence of alcohol-related tweets with alcohol-related outcomes reported from the National Institute of Alcoholism and Alcohol Abuse, Centers for Disease Control Behavioral Risk Factor Surveillance System, county health rankings, and the National Industry Classification System.

**Methods:** The Bidirectional Encoder Representations From Transformers neural network learned to classify 1,302,524 Reddit posts as either alcohol-related or control subreddits. The trained model identified 24 alcohol-related hashtags from an unlabeled data set of 843,769 random tweets. Querying alcohol-related hashtags identified 25,558,846 alcohol-related tweets, including 790,544 location-specific (geotagged) tweets. We calculated the correlation between the prevalence of alcohol-related tweets and alcohol-related outcomes, controlling for confounding effects of age, sex, income, education, and self-reported race, as recorded by the 2013-2018 American Community Survey.

**Results:** Significant associations were observed: between alcohol-hashtagged tweets and alcohol consumption ($P=.01$) and heavy drinking ($P=.005$) but not binge drinking ($P=.37$), self-reported at the metropolitan-micropolitan statistical area level; between alcohol-hashtagged tweets and self-reported excessive drinking behavior ($P=.03$) but not motor vehicle fatalities involving alcohol ($P=.21$); between alcohol-hashtagged tweets and the number of breweries ($P<.001$), wineries ($P<.001$), and beer, wine, and liquor stores ($P<.001$) but not drinking places ($P=.23$), per capita at the US county and county-equivalent level; and between alcohol-hashtagged tweets and all gallons of ethanol consumed ($P<.001$), as well as ethanol consumed from wine ($P<.001$) and liquor ($P=.01$) sources but not beer ($P=.63$), at the US state level.

**Conclusions:** Here, we present a novel natural language processing pipeline developed using Reddit's alcohol-related subreddits that identify highly specific alcohol-related Twitter hashtags. The prevalence of identified hashtags contains interpretable information about alcohol consumption at both coarse (eg, US state) and fine-grained (eg, metropolitan-micropolitan statistical area level and county) geographical designations. This approach can expand research and deep learning interventions on alcohol abuse and other behavioral health outcomes.

## Introduction

### Background

Alcohol-related causes are the third leading preventable cause of death in the United States, and alcohol abuse contributes to many adverse health outcomes, particularly on the developing brain [1-4]. The rise of alcohol-related content on Twitter is alarming, with over half of young adults participating in a study [5] posting alcohol-related content. Social media use and alcohol consumption are common behaviors; the prevalence rates of Twitter, Reddit, and annual alcohol use for US adults are 22%, 11%, and 70%, respectively [6,7]. Internet- and social media–based interventions are scalable and efficient approaches for developing practical tools for treating and monitoring alcohol abuse, especially for at-risk adolescents and young adults [8-14]. However, identifying high-risk areas for efficient and helpful monitoring along with population-level interventions remains a difficult task, in part because of survey bias [15-17].

Text-based *hashtags* are common among many popular social media platforms such as Twitter, Instagram, and TikTok. Individuals use hashtags to categorize, label, organize, and discover posts and content [18]. Previous studies have indicated that study-specific hashtags are useful for mental health research [19]. For example, sexual abuse and harassment (#MeToo), breast cancer (#breastcancer), HIV (#HIV), miscarriages (#ihadamiscarriage), tobacco use (#Vapelife), and viral pandemics (#COVID-19) are some of the many important health outcomes that have been previously studied using hashtags on Twitter [20-28]. Other social media platforms such as Reddit contain specific *themed* communities where interested users discuss a particular topic. In contrast to hashtags, themed communities on websites such as Reddit represent posts related to exactly 1 topic of interest. Like hashtags, these communities, such as *r/cripplingalcoholism*, *r/depression*, or *r/opiates* Reddit subreddits and *HIV* Baidu Tieba bar, contain information that can target and understand behavioral health and disease [29-33]. In addition to hashtags and subreddits, some social media platforms allow for *geotagging* or sharing a user's geographical latitude and longitude coordinates in a post. Geotags have been used in social media research to identify geographically relevant information from social media data [34-36].

### Previous Work

Although prior studies have identified specific hashtags or themed communities for studying behavioral health outcomes, many insights are platform-specific. Although helpful information regarding a behavior of interest or themed community may be available on one platform, there may not be such knowledge available on a different platform. Many previous methods examining alcohol content on social media use data from a single platform [5,37-42]. Single-platform analyses may limit discoveries and interventions to only a fraction of the population at risk. There is a growing need for behavioral health researchers working with social media data to incorporate analyses from many sources [43,44]. Although some studies have examined alcohol content on multiple platforms, many methods need survey data from known active users from each source or additional manual annotation [45-47]. The ability and insights gained from using deep learning methods to learn from a large number of posts from specific communities (ie, Reddit subreddits) to predict alcohol-related content on a different platform (ie, Twitter) remain unclear.

Many previous studies that identified alcohol-related language on social media platforms relied on training on extrinsic labels, such as survey responses. Reliance on self-report data is problematic as alcohol consumption is subject to bias, particularly among the youth [15,16,48]. In addition, approaches that use an outcome of interest to both train and evaluate a model (eg, identifying and evaluating alcohol-related hashtags or keywords based on enrichment in regions with higher self-reported alcohol content) may not be generalizable to other related outcomes [49].

Other approaches for studying alcohol content involve identifying a sample as being alcohol-related based on the identification of keywords. Keyword approaches have distinct benefits, such as interpretability. However, identifying text from keywords may rely on standard and predefined terms (eg, searching *drunk*), training on self-report data, or manual review [37,42,49-51]. Classification of social media posts based on previously defined keywords or vector representations (eg, Word2vec) is not as useful when the average length of sequences is small and has out-of-training vocabulary [52-54]. Training on nonspecific platform information alone may fail to capture relevant keywords, especially for rarer outcomes not prominent in the heterogeneity of random and unlabeled social media chatter [55]. In addition, predefined keywords or word vectors may fail to capture slang or the different language structures between Reddit and Twitter [56].

One recent contribution in natural language processing (NLP) is the Bidirectional Encoder Representations From Transformers (BERT) neural network, which has demonstrated superior performance on a wide variety of social media NLP tasks [57-61]. BERT focuses on learning by analyzing sentences with randomly masked words. This masked language model deconstructs larger strings into smaller tokens and is ideal for dealing with hashtags and other platform-unique token structures [57]. Before developing BERT, previous models, such as long-short–term memory networks, logistic regression, Word2vec similarity, and latent Dirichlet allocation, were not well suited to process unknown words and hashtag structures. For example, some previous NLP studies on social media either removed hashtags, represented them as universal tokens, or removed # from strings, with no importance given to hashtags (eg, *#ilovebeer* represented as " " (space), *HASHTAG*, or *ilovebeer*, respectively) [62-64]. In contrast, using hashtags and themed communities as explicit labels in a deep learning

XSL•FO

**RenderX**

architecture allows for identifying relevant, platform-specific hashtags that can identify posts that indicate the behavior of interest. In addition, the use of these structures adds a layer of interpretability to our trained neural networks, which are commonly criticized as noninterpretable black boxes [65].

Other previous social media text mining methods implementing deep learning often involve training platform-specific models. One issue with this approach is that each platform's training models require an extensive amount of usually labeled data from that platform [66-72]. In addition, although deep learning models have been successful at many tasks, training platform-specific deep networks such as BERT (containing >100 million parameters) is extremely energy- and cost-intensive, and $CO_2$ emissions from training BERT models have raised concerns about their environmental impact [73]. Optimal methods for translating information from previously trained social media deep learning models to discern insights from separate social media platforms remain a relatively unexplored research area.

### The Goal of This Study

We aim to examine the effectiveness of using thematic structures in a deep learning framework to identify alcohol-related behaviors across different social media platforms. First, we trained on Reddit subreddits to identify alcohol-related targets on another social media platform (Twitter) with a different thematic structure (hashtags). Next, we determined whether the hashtags predicted by the model correlate to known alcohol-related outcomes, including self-reported drinking status, alcohol outlet density, and estimated gallons of ethanol consumed, after controlling for confounding effects of age, sex, income, education, and self-reported race. We show that these data-driven hashtags contain interpretable information about alcohol consumption in the United States. Finally, we present validated and queryable hashtags from our model that behavioral health researchers can use as a starting point for the identification of alcohol-related content on Twitter, Reddit, and other social media platforms.

## Methods

### Overview of the NLP Pipeline

This study fine-tuned a BERT neural network as a binary classifier to predict Reddit post titles as belonging to either alcohol-related communities or a random subreddit. Next, we applied the Reddit-trained network to a smaller set of random, unlabeled Twitter posts to identify 24 hashtags that were significantly associated with alcohol content. We identified 25,558,846 tweets that contained at least one alcohol-related hashtag for the period between 2010 and 2019. A total of 1,412,041 alcohol-related tweets included latitude and longitude data from *geotagging*. The locations of 790,544 geotagged tweets from 2929 US counties and county equivalents were identified using data from the 2017 US Census Shapefiles database [74,75]. Finally, we examined the relationship between the prevalence of alcohol-related tweets per population and various outcome measures related to alcohol consumption, including self-reported alcohol consumption and alcohol outlet density. Figure 1 demonstrates an overview of our NLP pipeline. Figure 2 illustrates the choropleth of population-normalized alcohol-hashtagged tweets for US states and Washington, DC.

**Figure 1.** Overview of the methodological pipeline. A bidirectional encoder representation from transformers model trained to classify posts as either 18 alcohol-related or control subreddits. The bidirectional encoder representations from transformers model was applied to a set of tweets containing at least one hashtag. The prediction results were analyzed to find 24 significantly enriched hashtags as positive predictions (ie, prediction probability ≥0.5). Tweets posted between 2010 and 2020 with an alcohol-related hashtag were collected and filtered on geotagged location. BERT: Bidirectional Encoder Representations From Transformers.
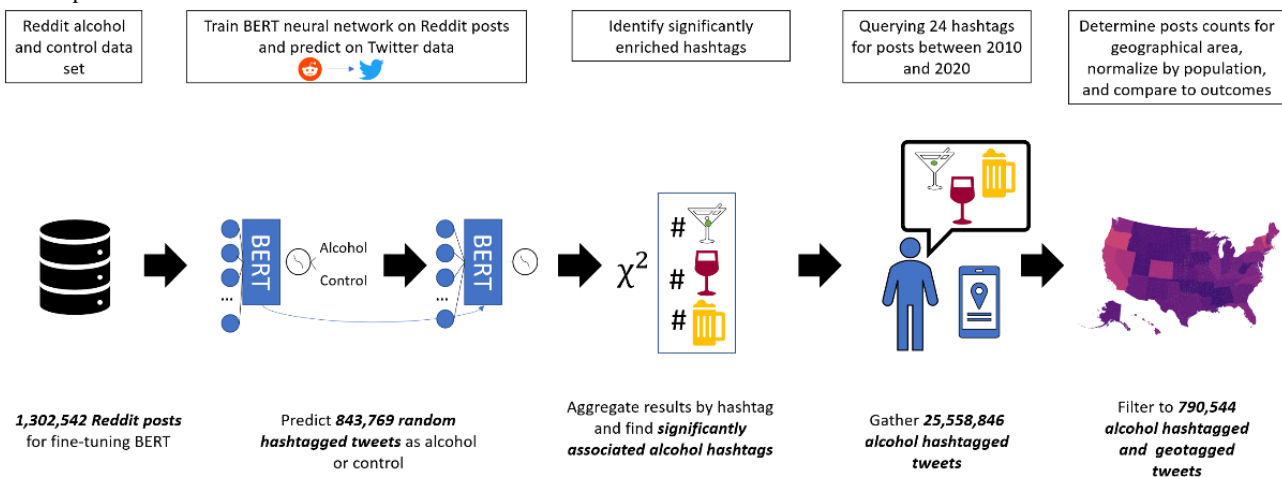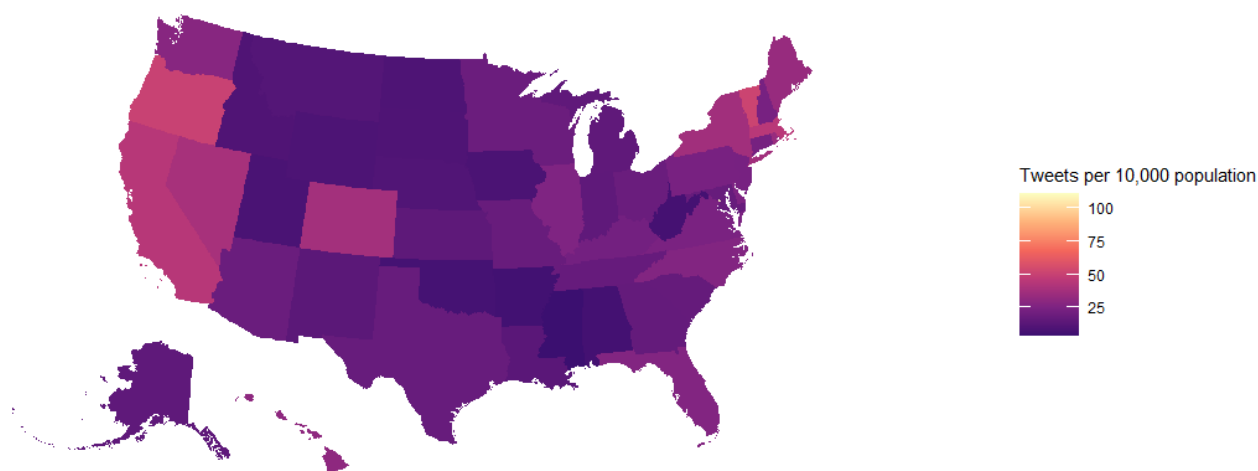
**Figure 2.** Choropleth of the US state and Washington, DC tweets with alcohol-related hashtags per 10,000 persons.



## Reddit Data Set and BERT Training

A large amount of alcohol-related training data were extracted from Reddit subreddits with the pushshift application programming interface (API) previously used in social media research [63,76]. On a subreddit, community moderators create description texts that contain links to other, usually related, subreddits. Scraping the description pages for all subreddits containing at least 1000 posts for any links to *r/drunk*, one of the most popular alcohol-related subreddits, and all links from *r/drunk* to other subreddits yielded 17 alcohol-related subreddits. A total of 651,271 post titles from the following 18 subreddits were used as positive alcohol labels for model training: *r/cripplingalcoholism*, *r/vodka*, *r/oldtimehockey*, *r/alcohol*, *r/beer*, *r/bourbon*, *r/homebrewing*, *r/drinkinggames*, *r/wine*, *r/beercirclejerk*, *r/gin*, *r/scotch*, *r/liquor*, *r/showerbeer*, *r/absinthe*, *r/firewater*, *r/beercanada*, and *r/drunk*.

Negative alcohol (control) posts were gathered by querying 651,271 random posts posted in all other subreddits, excluding the 18 alcohol-related subreddits. Training 79.99% (521,016/651,271), validation 9.99% (65,127/651,271), and testing 9.99% (65,127/651,271), data sets were generated for developing and evaluating the model—a binary classifier trained for posts belonging to either alcohol-related subreddits or other random subreddits. The model fine-tuned a pretrained BERT model with 12 layers and 768 hidden units in PyTorch on an NVIDIA TITAN Xp graphics processing unit using a batch size of 64 for approximately 5 weeks [77].

## Twitter Data Set and Identification of Hashtags

The Twitter API provides tweet information from 7 days before a query. Randomly selected tokens in the Twitter GLoVE word embedding dictionary and their respective hashtags (ie, a string that starts with #) were queried using the Twitter API to identify recently posted tweets containing that word or hashtag [78]. Each identified hashtag in the data set was requeried to ensure that it was monitored for at least 2 weeks. The initial *random* Twitter data set comprised 843,769 random hashtag-containing tweets posted between January 2019 and October 2019. The Reddit-trained BERT model was applied to this data set to obtain binary predictions for each tweet. A chi-square test identified 24 significant alcohol-related hashtags from posts predicted to be alcohol-positive (ie, final softmax layer prediction $P$ value of ≥0.5) relative to posts predicted as negative (ie, final softmax prediction $P$ value of <0.5) using a one-tailed *greater* test. We included only hashtags with 5 or more occurrences and applied the Benjamini-Hochberg algorithm for multiple hypothesis correction using a 0.05 false discovery rate, a common approach for multiple hypothesis corrections in social media data analyses [37,42,79-85]. The analysis resulted in 24 hashtags, as indicated in Textbox 1. GetOldTweets, a Python package widely used in social media research, was used to identify 25,558,846 alcohol-hashtagged tweets posted throughout 10 years (between 2010 and 2020) containing at least one significant alcohol-related hashtag [55,86,87].

**Textbox 1.** Alcohol-related hashtags extracted by our Reddit-trained classifier according to alcohol category.

---

**Beer hashtags**

- *craftbeer*, *beer*, *ncbeer*, *brewery*, *stout*, *beeroclock*, *beergeek*, *beerporn*, *beers*, *instabeer*, *beertime*, *beerstagram*, *beerlover*, *beersnob*

**Wine hashtags**

- *winetasting*, *wine*, *winelover*, *wines*, *redwine*

**Liquor hashtags**

- *bourbon*, *whiskey*, *whisky*

**Multiple or ambiguous hashtags**

- *drinklocal*, *drunktwitter*

---

## Geographical Identification of the Prevalence of Alcohol-Related Hashtags

Next, we tested whether the knowledge of 24 significant alcohol hashtags could uncover information on alcohol-related outcomes in the United States. Alcohol-hashtagged tweets were filtered to 790,544 *geotagged* tweets containing longitude and latitude coordinate locations and mapped to metropolitan-micropolitan statistical areas (MMSAs), US county and county equivalents, and US states and Washington, DC. The total number of alcohol-hashtagged tweets in an area divided by the mean of the population estimates from the 2013-2018 American Community Survey yielded population-normalized alcohol-related hashtag prevalence.

We then tested the association between geographical prevalence of alcohol-related hashtags and alcohol outcomes. Spearman rho, a ranked nonparametric measure that is more robust to outliers than Pearson correlation, is used to report crude (nonadjusted) correlations [88]. Potential confounding variables previously studied in alcohol and social media use include race and sex distribution, median age, education, and income [37,89]. A linear regression analysis evaluated the relationship between the number of tweets per population and alcohol-related outcomes after including terms to control for confounding effects. Specific confounding variables from the 2013-2018 5-year American Community Survey report included *Percent Reporting White*, *Percent Reporting Black*, *Percent Reporting Hispanic*, *Median Income*, *Percent High School Education*, *Percent Bachelor's Degree Education*, and *Males/100 Females* [90]. All alcohol outcomes and confounding variables represented the most recent estimation of alcohol consumption and related behavior at the time of this study.

## Metropolitan-Micropolitan Statistical Areas

MMSAs are US Census Bureau designations of concentrated urban centers that may be the integrated areas of multiple cities and states (eg, the single *Washington-Arlington-Alexandria, DC-VA-MD-WV MSA* contains 3 US states and Washington, DC) [74]. The Behavioral Risk Factor Surveillance System publishes reports of survey responses at selected MMSAs for the following categories [91]:

- any alcohol consumption, defined as at least one alcoholic drink in the last 30 days;

- binge drinking behavior, defined as drinking >5 drinks in 1 event for men or >4 drinks in 1 event for women;
- heavy drinking, defined as drinking >1 drink per day for women or >2 drinks per day per man;

All yearly records from 2010-2019 for each MMSA were averaged to obtain a single number for outcome measurements.

## US County and County Equivalent

Primary US county outcomes were gathered from the University of Wisconsin Population Health Institute County Health Rankings and Roadmaps 2020 data, which included the estimates of excessive drinking, defined as the percentage reporting either binge or heavy drinking behavior as well as measurements of the percentage of motor vehicle fatalities that involved alcohol for the period between 2013 and 2018 [92]. In addition, data from the North American Industry Classification System provided by 2017 County Business Patterns (US Census) was used for the number of *Drinking Places (Alcohol Beverages)*; *Wineries*; *Breweries*; and *Beer*, *Wine*, *and Liquor stores* (North American Industry Classification System codes 722410, 312130, 312120, 445310, respectively) present in each county [93]. Counties were included if they contained at least one tweet and an average reported population >1000 between 2013 and 2018.

## US States and Washington DC

Twitter posts containing alcohol-related hashtags were aggregated by state and compared with the National Institute on Alcohol Abuse and Alcoholism's 2018 report, *Apparent Per Capita Alcohol Consumption: National, State, and Regional Trends*. This report predicts gallons of ethanol consumption based on alcohol sales and taxation data, separated for the consumption of wine, beer, or liquor products [94]. To determine which hashtags may be useful for detecting individual preferences of alcohol consumption, we calculated the correlation between the consumption of alcohol from different sources of alcoholic drinks (beer, wine, and liquor) and the prevalence of 19 beer, 5 wine, and 3 liquor-specific hashtags, as indicated in Textbox 1.

## *Results*

Table 1 demonstrates the results from the analysis of alcohol-related hashtags and alcohol-related outcomes. The
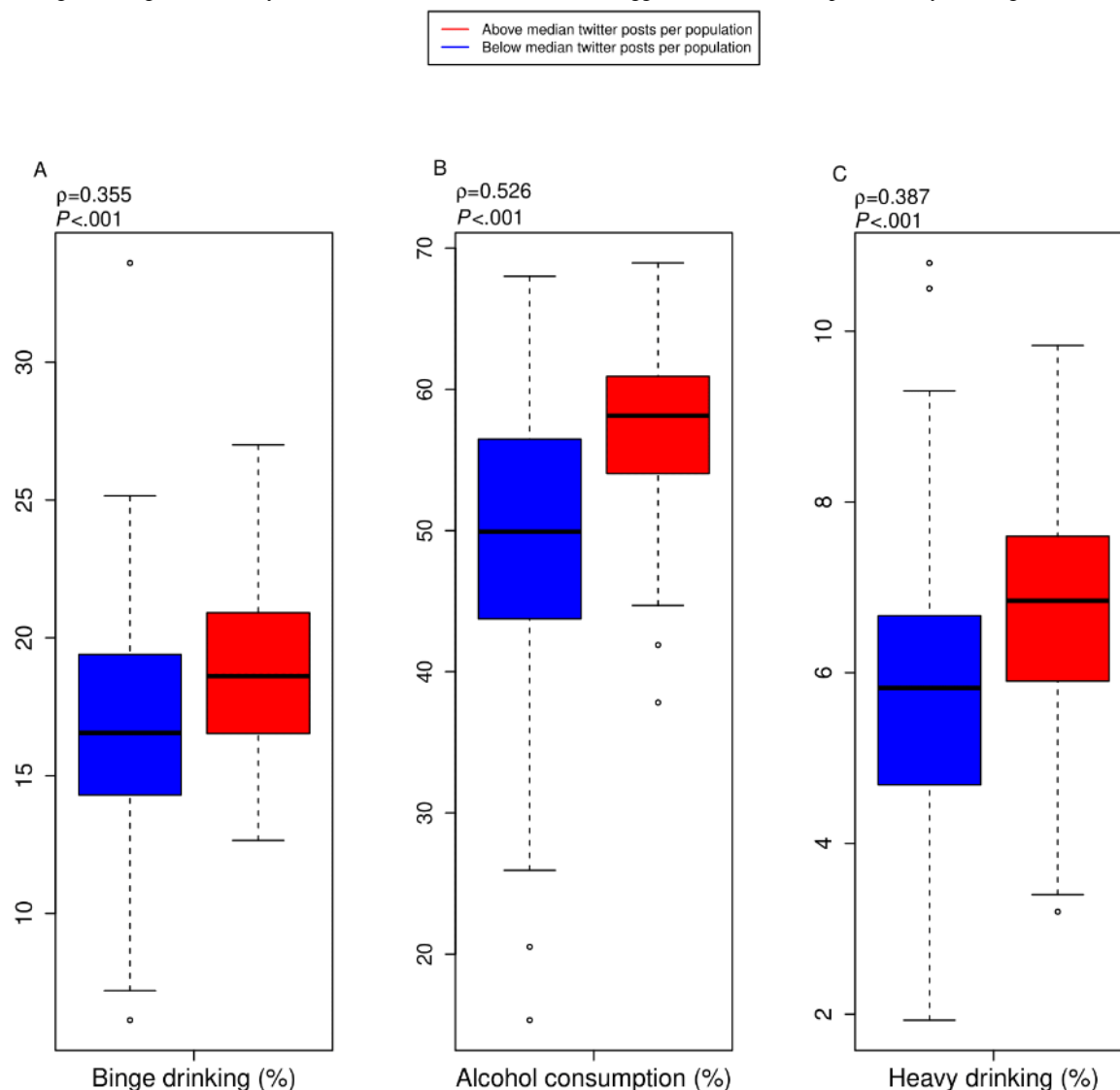
XSL•FO

**RenderX**

number of geotagged and alcohol-hashtagged tweets per population significantly correlated with many alcohol-related outcomes, including self-reported measures of individuals (1) reporting any alcohol consumption within 30 days ($P<.001$), (2) meeting the criteria for heavy drinking ($P<.001$), and (3) meeting the criteria for binge drinking ($P<.001$) at the MMSA level (Figure 3). However, the relationship between MMSA tweets and binge drinking level was not significant after adjusting for confounding effects ($P=.37$).

**Table 1.** Spearman correlation and linear regression results between the number of tweets per population and alcohol-related behavior and health indicators.

| Outcome | Spearman correlation | | Adjusted regression | | Sample size, n |
|---|---|---|---|---|---|
| | ρ | *P* value | Coefficient β | *P* value | |
| **Metropolitan-micropolitan statistical area** | | | | | |
| Alcohol consumption | 0.526 | <.001 | 1038 | .01 | 179 |
| Binge drinking | 0.355 | <.001 | 184.0 | .37 | 179 |
| Heavy drinking | 0.387 | <.001 | 244.8 | .005 | 179 |
| **County and equivalent** | | | | | |
| Excessive drinking | 0.377 | <.001 | 32.8 | .03 | 2641 |
| Percentage of alcohol motor vehicle fatality | 0.063 | .002 | 110.0 | .21 | 2641 |
| Drinking places (alcoholic beverages) per capita | −0.177 | <.001 | −2.18e–03 | .23 | 1479 |
| Breweries per capita | 0.263 | <.001 | 1.86e–03 | <.001 | 334 |
| Wineries per capita | 0.130 | .05 | 2.73e–02 | <.001 | 228 |
| Beer, wine, and liquor stores per capita | −0.043 | .11 | 0.0039 | <.001 | 1444 |
| **US states and Washington, DC, all hashtags** | | | | | |
| Wine, gallons of ethanol per capita | 0.756 | <.001 | 74.11 | <.001 | 51 |
| Beer, gallons of ethanol per capita | −0.050 | .73 | 9.911 | .63 | 51 |
| Liquor, gallons of ethanol per capita | 0.320 | .01 | 62.54 | .03 | 51 |
| All sources, gallons of ethanol per capita | 0.437 | <.001 | 146.6 | <.001 | 51 |
| **US states and Washington, DC, hashtags stratified by alcohol category** | | | | | |
| Wine, gallons of ethanol per capita (5 hashtags) | 0.754 | <.001 | 214.6 | <.001 | 51 |
| Beer, gallons of ethanol per capita (19 hashtags) | −0.001 | .99 | 16.05 | .63 | 51 |
| Liquor, gallons of ethanol per capita (3 hashtags) | 0.140 | .33 | 338.0 | .01 | 51 |

**Figure 3.** Metropolitan-micropolitan statistical area correlations for alcohol-hashtagged tweets and percent self-reported alcohol consumption. (A) Number of alcohol-hashtagged tweets and self-reported alcohol consumption within 30 days (N=179); (B) number of alcohol-hashtagged tweets and self-reported binge drinking within 30 days (N=179); (C) number of alcohol-hashtagged tweets and self-reported heavy drinking within 30 days (N=179).



There was a significant correlation between the percentage of motor vehicle deaths reported as involving alcohol ($P<.001$) and aggregated measures of excessive drinking behavior ($P<.001$) at the county level (Figure 4). However, the relationship between alcohol-related motor vehicle fatalities and county tweets per population was not significant after adjusting for confounding effects ($P=.21$).

**Figure 4.** County- and county-equivalent–level correlations for alcohol-hashtagged tweets, alcohol motor vehicle fatalities, and self-reported excessive alcohol consumption. Counties are included if they contain at least one tweet and 1000 people (N=2641). (A) County-level correlations for alcohol-hashtagged tweets per population and percentage of motor vehicle fatalities involving alcohol. (B) County-level correlations for alcohol-hashtagged tweets per population and percent self-reporting excessive drinking (reporting either more than five alcoholic drinks on a single occasion for men or more than four alcoholic drinks on a single occasion for women, or more than two drinks per day for men or more than one drink per day for women).



There was a significant correlation between the number of alcohol-hashtagged tweets and wineries (*P*=.05), breweries (*P*<.001), and drinking places (alcoholic beverages; *P*<.001) but not beer, wine, and liquor stores (*P*=.11) per capita at the county level (Figure 5). However, after adjusting for confounding effects, there was a significant association between alcohol-related tweets per population and beer, wine, and liquor stores (*P*<.001) but not drinking places (*P*=.23) per capita.

**Figure 5.** County- and county-equivalent–level correlations for alcohol-hashtagged tweets per person and alcohol-serving outlets, as reported by the North American Industry Classification System. Counties are included if they have at least one (1) tweet in our data set, one (1) alcohol outlet, and contain a population of 1000. (A) Wineries per 10,000 people (n=228 counties); (B) breweries per 10,000 people (n=334); (C) liquor stores per 10,000 people (n=1444); (D) drinking places (alcoholic beverages) per 10,000 people (n=1479).



There was a significant correlation between the prevalence of alcohol-hashtagged tweets and gallons of wine (*P*<.001), liquor (*P*=.01), and overall gallons of consumption (*P*<.001) at the state level (Figure 6). However, the association with all alcohol-hashtagged tweets and gallons of beer consumed was not significant (*P*=.63).

The prevalence of five wine hashtags had a significant association with wine consumption at the state level (*P*<.001), but 3 liquor hashtags and 19 beer hashtags did not have a significant relationship with liquor (*P*=.33) and beer (*P*=.99) consumption at the state level. However, there was a significant relationship between gallons of liquor consumed and the prevalence of 3 liquor hashtags (*P*=.01) after controlling for confounding effects (Figure 6).

**Figure 6.** Gallons of ethanol consumed for US states and alcohol-hashtagged tweets. Population normalization was performed using the average reported population between 2010 and 2018, as reported by the National Institute of Alcohol Abuse and Alcoholism. Top, all 24 hashtags used; gallons of ethanol consumed from (A) wine; (B) beer; (C) liquor; and (D) all alcohol sources and Twitter posts per population. Bottom, only specific hashtags, gallons of ethanol from wine, beer, and liquor, and Twitter posts per population from specific hashtags: (E) wine (5 wine hashtags); (F) beer (14 beer hashtags); (G) liquor (3 liquor hashtags).



## Discussion

### Principal Findings

We demonstrated that information from alcohol-related subreddits could identify alcohol-related hashtags that correlated with multiple alcohol-related outcomes. The prevalence of geotagged tweets containing these significantly alcohol-related hashtags correlates with alcohol-related behaviors and alcohol outlet density, which are associated with adverse outcomes, including deaths from motor vehicle crashes and excessive drinking [95-97]. This approach has distinct benefits for studying

alcohol-related outcomes. Compared with other approaches to detect alcohol consumption on Twitter, the pipeline presented here is trained using relatively few interpretable input and output parameters, specifically, 18 subreddits and 24 hashtags. The NLP pipeline identified language associated with alcohol abuse on Twitter without manual annotation or predefined keywords and with only the knowledge of relevant subreddits. These results indicate that alcohol-related language, when defined by inclusion into the themed communities of subreddits or containing alcohol-related hashtags, can be used to understand population-level behavior in multiple geographic areas with different population granularity.

Qualitatively, the model presented here detects a wide variety of different alcohol-related hashtags, including slang. Tweets containing these hashtags capture a large amount of information regarding alcohol consumption behavior at a broad population level. The set of hashtags resulting from this study is useful for future alcohol-related research on Twitter and identifying relevant hashtags or community forum labels on other platforms.

A notable benefit of using hashtags and subreddits as platform-specific labels for studying alcohol-related outcomes is interpretability. Although deep learning models excel at understanding massive quantities of data, many NLP and deep learning models rely on complex feature representations to classify or characterize text. However, this approach is not easily interpretable [65,98]. We have demonstrated that by treating subreddits and hashtags as learnable labels, it is possible to directly use hashtags as interpretable features in our NLP pipeline for social media data to understand alcoholic beverage preferences while still learning from a large amount of data.

Consumption of different alcohol types is associated with a variety of both beneficial and detrimental outcomes [99-101]. For example, wine consumption is associated with protection from cardiovascular diseases; however, the cause of this effect may be dietary factors and other lifestyle choices [99,102]. Other studies have associated preferential beer and liquor consumption with adverse outcomes, such as dangerous drinking and other risky behaviors [101,103,104]. Notably, although many studies have examined overall alcohol mentions on Twitter, few models created have been explicitly examined in terms of differences in the types of alcoholic beverages mentioned.

Our study indicates that information capturing consumption of wine and liquor is directly observable using social media data, as shown by the significant associations between the prevalence of 5 wine-related hashtags and the amount of wine consumed, as well as the number of posts containing at least 1 of the 3 liquor-related hashtags and liquor consumed. However, there was no significant relationship between beer consumption per capita and the number of alcohol-hashtagged Twitter posts in an area. The results here indicate that our model can detect certain types of alcohol consumption behavior (wine and liquor consumption) on Twitter using the interpretability of hashtags but not others (beer consumption). It remains unclear whether our results indicate a bias in our model's methodological choice (eg, use of hashtags or training procedure) or a difference in social media populations that prefer different alcohol beverages.

The difference in correlations might be because of several variables, including the existence of confounding factors related to the prevalence of social media use in the underlying populations that preferentially consume beer over other alcohol types, differences in perceived acceptance of beer consumption behavior, or because of other factors that may confound alcohol and social media posting [89]. Alcohol preference is an example of having interpretable hashtag representations for a given model that may help identify behavioral differences associated with an outcome of interest. This evidence suggests that similar models trained on Twitter may detect alcohol wine and liquor consumption but not beer.

## Comparison With Previous Works

Many previous studies have used hashtags as target labels. However, they mostly rely on a predefined set of hashtags that may not be data-driven or require extensive expert annotation instead of taking advantage of topic-specific sources and social media content on other public social media platforms [40,105-107]. Notably, using a predefined number of hashtags could be biased and too narrow for capturing relevant information, potentially missing informative hashtags for exposure or outcome. Traditional keyword approaches may fail to capture various pieces of information from slang and novel hashtags from platform-specific languages as they are created and popularized. In addition, keyword databases may not exist for all outcomes of interest. Semantic similarity measures, such as Word2vec, may identify hashtags with similar contexts; however, integrating vector representations may lose valuable information relative to training over individual samples, and prediction probabilities or certainty are not readily observable. In contrast, this study indicates that a model trained on a large set of data relevant to the behavior of interest and its application to an unlabeled data set from a different platform can identify data-driven hashtags related to that behavior. This ability to learn hashtags from data is critical, as new hashtags are created every day and may differ substantially between platforms.

Many social media platforms are directly searchable using hashtags, allowing the ability to gather many highly specific posts instead of gathering a large number of nonspecific posts to identify relevant hashtags, keywords, or alcohol content based on available prevalence data. Although the latter approach has shown success in studying alcohol-related behavior previously, methods to extend the analysis to alternative but potentially related outcomes, platforms, or different geographical designations remain unclear [39,40,108]. The generalization of such models necessitates the creation of individual models for each particular outcome and geographic area. Finally, these models may fail to take advantage of extensive research on themed communities to understand alcohol use and other outcomes of interest [30,39,40,76,106,109,110]. In contrast, the method outlined here may help create more efficient public health interventions to analyze the alcohol consumption behavior for a given geographic area of interest, such as a city or hospital catchment area. This approach is particularly useful when the relevant language is dynamic or contains area-specific slang, making previously established dictionary-based methods incomplete or impractical. In particular, our proposed methodology for identifying hashtags using a previously trained

deep learning model can be useful for detecting alcohol consumption behavior on various social media platforms.

## Limitations and Future Work

There are some limitations to the model and approach proposed in this study. Other health-related behaviors that are not discussed frequently on social media may be more difficult to ascertain and translate between platforms using hashtags. Any model trained on a limited number of social media platforms may be confounded by differences in user preferences, such as age or socioeconomic status. Furthermore, this method relies on identifying known or previously studied subreddits, which may not be suitable for outcomes without known relevant subreddits. Furthermore, we did not compare our models' performance with the BERT large model or other deep learning alternatives, which can have a different performance for our task. In addition, including additional social media data in model development and using domain knowledge from ontologies, controlled vocabularies, lexicons, and relevant rules and regular expressions can further improve the presented results.

As future work, we plan to extend this study to other mental and behavioral health topics, such as depression and substance use, and other social media platforms that use hashtags, such as Facebook and Instagram. Deep learning has previously been used to combine the analysis of images and texts from social media users [38]. In our future work, we will expand the presented architecture to include other data modalities, such as images and videos, to increase screening capabilities.

## Conclusions

These results indicate that using alcohol-related subreddits as learnable labels to train a BERT neural network can capture interpretable, alcohol-related language on Twitter. Our study suggests a significant correlation between the prevalence of alcohol-related geotagged Twitter hashtags and alcohol-related behaviors as measured by self-reported alcohol consumption, alcohol preferences, and alcohol outlet prevalence. This method has the unique advantages of previous methods, including allowing examination at the MMSA, US county, and US state level for different alcohol-related outcomes. These results suggest that using previously studied hashtags and subreddits as learnable targets in a machine learning framework could expand public health outreach efforts and epidemiology research, particularly for monitoring behavior related to alcohol consumption.

## Authors' Contributions

Study concept and design: BJR and SH; analysis and interpretation of data: BJR and SH; collection or assembly of data: BJR; drafting of the manuscript: BJR and SH; critical revision of the manuscript for important intellectual content: BJR and SH; funding: BJR and SH; administrative, technical, and logistic support: SH; and study supervision: SH.

## Conflicts of Interest

None declared.

## References

1.  Osna NA, Donohue TM, Kharbanda KK. Alcoholic liver disease: pathogenesis and current management. Alcohol Res 2017;38(2):147-161 [FREE Full text] [Medline: 28988570]

2.  O'Keefe JH, Bhatti SK, Bajwa A, DiNicolantonio JJ, Lavie CJ. Alcohol and cardiovascular health: the dose makes the poison…or the remedy. Mayo Clin Proc 2014 Mar;89(3):382-393. [doi: 10.1016/j.mayocp.2013.11.005] [Medline: 24582196]

3.  Alcohol and public health: Alcohol-Related Disease Impact (ARDI). Centers for Disease Control and Prevention (CDC). URL: https://nccd.cdc.gov/DPH_ARDI/default/default.aspx [accessed 2021-08-24]

4.  Zhao Q, Sullivan EV, Honnorat N, Adeli E, Podhajsky S, De Bellis MD, et al. Association of heavy drinking with deviant fiber tract development in frontal brain systems in adolescents. JAMA Psychiatry 2021 Apr 01;78(4):407-415. [doi: 10.1001/jamapsychiatry.2020.4064] [Medline: 33377940]

5.  Litt DM, Lewis MA, Spiro ES, Aulck L, Waldron KA, Head-Corliss MK, et al. #drunktwitter: examining the relations between alcohol-related Twitter content and alcohol willingness and use among underage young adults. Drug Alcohol Depend 2018 Dec 01;193:75-82 [FREE Full text] [doi: 10.1016/j.drugalcdep.2018.08.021] [Medline: 30343237]

6.  Perrin A, Anderson M. Share of U.S. adults using social media, including Facebook, is mostly unchanged since 2018. Pew Research Center. URL: https://www.pewresearch.org/fact-tank/2019/04/10/share-of-u-s-adults-using-social-media-including-facebook-is-mostly-unchanged-since-2018/ [accessed 2021-08-24]

7.  2018 National Survey of Drug Use and Health (NSDUH) releases. SAMHSA. URL: https://www.samhsa.gov/data/release/2018-national-survey-drug-use-and-health-nsduh-releases [accessed 2021-08-24]

8.  Bonar EE, Schneeberger DM, Bourque C, Bauermeister JA, Young SD, Blow FC, et al. Social media interventions for risky drinking among adolescents and emerging adults: protocol for a randomized controlled trial. JMIR Res Protoc 2020 May 13;9(5):e16688 [FREE Full text] [doi: 10.2196/16688] [Medline: 32401225]

XSL•FO

RenderX

9.    Arigo D, Pagoto S, Carter-Harris L, Lillie SE, Nebeker C. Using social media for health research: methodological and ethical considerations for recruitment and intervention delivery. Digit Health 2018 May 7;4:2055207618771757 [FREE Full text] [doi: 10.1177/2055207618771757] [Medline: 29942634]

10.   Eysenbach G. Infodemiology: the epidemiology of (mis)information. Am J Med 2002 Dec 15;113(9):763-765. [doi: 10.1016/s0002-9343(02)01473-0] [Medline: 12517369]

11.   Mavragani A. Infodemiology and infoveillance: scoping review. J Med Internet Res 2020 Apr 28;22(4):e16206 [FREE Full text] [doi: 10.2196/16206] [Medline: 32310818]

12.   Eysenbach G. Infodemiology and infoveillance: framework for an emerging set of public health informatics methods to analyze search, communication and publication behavior on the internet. J Med Internet Res 2009 Mar 27;11(1):e11 [FREE Full text] [doi: 10.2196/jmir.1157] [Medline: 19329408]

13.   Martino F, Brooks R, Browne J, Carah N, Zorbas C, Corben K, et al. The nature and extent of online marketing by big food and big alcohol during the COVID-19 pandemic in Australia: content analysis study. JMIR Public Health Surveill 2021 Mar 12;7(3):e25202. [doi: 10.2196/25202] [Medline: 33709935]

14.   Song T, Qian S, Yu P. Mobile health interventions for self-control of unhealthy alcohol use: systematic review. JMIR Mhealth Uhealth 2019 Jan 29;7(1):e10899 [FREE Full text] [doi: 10.2196/10899] [Medline: 30694200]

15.   Dawson DA, Goldstein RB, Pickering RP, Grant BF. Nonresponse bias in survey estimates of alcohol consumption and its association with harm. J Stud Alcohol Drugs 2014 Jul;75(4):695-703 [FREE Full text] [doi: 10.15288/jsad.2014.75.695] [Medline: 24988268]

16.   Sartor CE, Bucholz KK, Nelson EC, Madden PA, Lynskey MT, Heath AC. Reporting bias in the association between age at first alcohol use and heavy episodic drinking. Alcohol Clin Exp Res 2011 Aug;35(8):1418-1425 [FREE Full text] [doi: 10.1111/j.1530-0277.2011.01477.x] [Medline: 21438885]

17.   Davis CG, Thake J, Vilhena N. Social desirability biases in self-reported alcohol consumption and harms. Addict Behav 2010 Apr;35(4):302-311. [doi: 10.1016/j.addbeh.2009.11.001] [Medline: 19932936]

18.   Zappavigna M. Searchable talk: the linguistic functions of hashtags. Soc Semiot 2015 Jan 09;25(3):274-291. [doi: 10.1080/10350330.2014.996948]

19.   Berry N, Lobban F, Belousov M, Emsley R, Nenadic G, Bucci S. #WhyWeTweetMH: understanding why people use twitter to discuss mental health problems. J Med Internet Res 2017 Apr 05;19(4):e107 [FREE Full text] [doi: 10.2196/jmir.6173] [Medline: 28381392]

20.   Modrek S, Chakalov B. The #MeToo Movement in the United States: text analysis of early Twitter conversations. J Med Internet Res 2019 Sep 03;21(9):e13837 [FREE Full text] [doi: 10.2196/13837] [Medline: 31482849]

21.   O'Neil A, Sojo V, Fileborn B, Scovelle AJ, Milner A. The #MeToo movement: an opportunity in public health? Lancet 2018 Jun 30;391(10140):2587-2589. [doi: 10.1016/S0140-6736(18)30991-7] [Medline: 30070210]

22.   Kero KM, Puuronen AH, Nyqvist L, Langén VL. Usability of two brief questions as a screening tool for domestic violence and effect of #MeToo on prevalence of self-reported violence. Eur J Obstet Gynecol Reprod Biol 2020 Dec;255:92-97 [FREE Full text] [doi: 10.1016/j.ejogrb.2020.10.024] [Medline: 33113404]

23.   Platt JR, Brady RR. #BCSM and #breastcancer: contemporary cancer-specific online social media communities. Breast J 2020 Apr;26(4):729-733. [doi: 10.1111/tbj.13576] [Medline: 31493301]

24.   Basch C, MacLean S. Breast cancer on Instagram: a descriptive study. Int J Prev Med 2019 Oct 9;10:166 [FREE Full text] [doi: 10.4103/ijpvm.IJPVM_36_19] [Medline: 32133084]

25.   George N, Britto D, Krishnan V, Dass L, Prasant H, Aravindhan V. Assessment of hashtag (#) campaigns aimed at health awareness in social media. J Educ Health Promot 2018 Sep 14;7:114 [FREE Full text] [doi: 10.4103/jehp.jehp_37_18] [Medline: 30271799]

26.   Mercier RJ, Senter K, Webster R, Henderson Riley A. Instagram users' experiences of miscarriage. Obstet Gynecol 2020 Jan;135(1):166-173. [doi: 10.1097/AOG.0000000000003621] [Medline: 31809440]

27.   Laestadius LI, Wahl MM, Cho YI. #Vapelife: an exploratory study of electronic cigarette use and promotion on instagram. Subst Use Misuse 2016 Oct 14;51(12):1669-1673. [doi: 10.1080/10826084.2016.1188958] [Medline: 27484191]

28.   Kudchadkar S, Carroll C. Using social media for rapid information dissemination in a pandemic: #PedsICU and coronavirus disease 2019. Pediatr Crit Care Med 2020 Aug;21(8):538-546 [FREE Full text] [doi: 10.1097/PCC.0000000000002474] [Medline: 32459792]

29.   Gkotsis G, Oellrich A, Velupillai S, Liakata M, Hubbard TJ, Dobson RJ, et al. Characterisation of mental health conditions in social media using informed deep learning. Sci Rep 2017 Mar 22;7:45141 [FREE Full text] [doi: 10.1038/srep45141] [Medline: 28327593]

30.   Tamersoy A, De Choudhury M, Chau D. Characterizing smoking and drinking abstinence from social media. HT ACM Conf Hypertext Soc Media 2015 Sep;2015:139-148 [FREE Full text] [doi: 10.1145/2700171.2791247] [Medline: 26640831]

31.   Liu C, Lu X. Analyzing hidden populations online: topic, emotion, and social network of HIV-related users in the largest Chinese online community. BMC Med Inform Decis Mak 2018 Jan 05;18(1):2 [FREE Full text] [doi: 10.1186/s12911-017-0579-1] [Medline: 29304788]

32.   Pandrekar S, Chen X, Gopalkrishna G, Srivastava A, Saltz M, Saltz J, et al. Social media based analysis of opioid epidemic using Reddit. AMIA Annu Symp Proc 2018 Dec 5;2018:867-876 [FREE Full text] [Medline: 30815129]

33. Kochan A, Ong S, Guler S, Johannson KA, Ryerson CJ, Goobie GC. Social media content of idiopathic pulmonary fibrosis groups and pages on Facebook: cross-sectional analysis. JMIR Public Health Surveill 2021 May 31;7(5):e24199 [FREE Full text] [doi: 10.2196/24199] [Medline: 34057425]

34. Gruebner O, Lowe SR, Sykora M, Shankardass K, Subramanian SV, Galea S. A novel surveillance approach for disaster mental health. PLoS One 2017 Jul 19;12(7):e0181233 [FREE Full text] [doi: 10.1371/journal.pone.0181233] [Medline: 28723959]

35. Jurdak R, Zhao K, Liu J, AbouJaoude M, Cameron M, Newth D. Understanding human mobility from Twitter. PLoS One 2015 Jul 8;10(7):e0131469 [FREE Full text] [doi: 10.1371/journal.pone.0131469] [Medline: 26154597]

36. Hussain A, Tahir A, Hussain Z, Sheikh Z, Gogate M, Dashtipour K, et al. Artificial intelligence-enabled analysis of public attitudes on Facebook and Twitter toward COVID-19 vaccines in the united kingdom and the United States: observational study. J Med Internet Res 2021 Apr 05;23(4):e26627 [FREE Full text] [doi: 10.2196/26627] [Medline: 33724919]

37. Curtis B, Giorgi S, Buffone AE, Ungar LH, Ashford RD, Hemmons J, et al. Can Twitter be used to predict county excessive alcohol consumption rates? PLoS One 2018 Apr 4;13(4):e0194290 [FREE Full text] [doi: 10.1371/journal.pone.0194290] [Medline: 29617408]

38. Hassanpour S, Tomita N, DeLise T, Crosier B, Marsch LA. Identifying substance use risk based on deep neural networks and Instagram social media data. Neuropsychopharmacology 2019 Feb;44(3):487-494 [FREE Full text] [doi: 10.1038/s41386-018-0247-x] [Medline: 30356094]

39. Marengo D, Azucar D, Giannotta F, Basile V, Settanni M. Exploring the association between problem drinking and language use on Facebook in young adults. Heliyon 2019 Oct 9;5(10):e02523 [FREE Full text] [doi: 10.1016/j.heliyon.2019.e02523] [Medline: 31667380]

40. Crocamo C, Viviani M, Bartoli F, Carrà G, Pasi G. Detecting binge drinking and alcohol-related risky behaviours from Twitter's users: an exploratory content- and topology-based analysis. Int J Environ Res Public Health 2020 Feb 26;17(5):1510 [FREE Full text] [doi: 10.3390/ijerph17051510] [Medline: 32111047]

41. Cavazos-Rehg PA, Krauss MJ, Sowles SJ, Bierut LJ. "Hey Everyone, I'm Drunk." An evaluation of drinking-related Twitter chatter. J Stud Alcohol Drugs 2015 Jul;76(4):635-643 [FREE Full text] [doi: 10.15288/jsad.2015.76.635] [Medline: 26098041]

42. Giorgi S, Yaden DB, Eichstaedt JC, Ashford RD, Buffone AE, Schwartz HA, et al. Cultural differences in Tweeting about drinking across the US. Int J Environ Res Public Health 2020 Feb 11;17(4):1125 [FREE Full text] [doi: 10.3390/ijerph17041125] [Medline: 32053866]

43. Stellefson M, Paige SR, Chaney BH, Chaney JD. Social media and health promotion. Int J Environ Res Public Health 2020 May 11;17(9):3323 [FREE Full text] [doi: 10.3390/ijerph17093323] [Medline: 32403215]

44. Sadah SA, Shahbazi M, Wiley MT, Hristidis V. A study of the demographics of web-based health-related social media users. J Med Internet Res 2015 Aug 06;17(8):e194 [FREE Full text] [doi: 10.2196/jmir.4308] [Medline: 26250986]

45. Hendriks H, Van den Putte B, Gebhardt WA, Moreno MA. Social drinking on social media: content analysis of the social aspects of alcohol-related posts on Facebook and Instagram. J Med Internet Res 2018 Jun 22;20(6):e226 [FREE Full text] [doi: 10.2196/jmir.9355] [Medline: 29934290]

46. Laws R, Hunt G, Antin TM. Social media platforms as a photo-elicitation tool in research on alcohol intoxication and gender. Nordisk Alkohol Nark 2018 Aug;35(4):288-303 [FREE Full text] [doi: 10.1177/1455072518781998] [Medline: 30245584]

47. Gupta H, Lam T, Pettigrew S, Tait RJ. The association between exposure to social media alcohol marketing and youth alcohol use behaviors in India and Australia. BMC Public Health 2018 Jun 13;18(1):726 [FREE Full text] [doi: 10.1186/s12889-018-5645-9] [Medline: 29895264]

48. Rehm J, Kilian C, Rovira P, Shield KD, Manthey J. The elusiveness of representativeness in general population surveys for alcohol. Drug Alcohol Rev 2021 Feb;40(2):161-165. [doi: 10.1111/dar.13148] [Medline: 32830351]

49. Nguyen QC, McCullough M, Meng H, Paul D, Li D, Kath S, et al. Geotagged US Tweets as predictors of county-level health outcomes, 2015-2016. Am J Public Health 2017 Nov;107(11):1776-1782. [doi: 10.2105/AJPH.2017.303993] [Medline: 28933925]

50. Krauss MJ, Grucza RA, Bierut LJ, Cavazos-Rehg PA. "Get drunk. Smoke weed. Have fun.": a content analysis of tweets about marijuana and alcohol. Am J Health Promot 2017 May;31(3):200-208 [FREE Full text] [doi: 10.4278/ajhp.150205-QUAL-708] [Medline: 26559715]

51. Allem J, Dormanesh A, Majmundar A, Rivera V, Chu M, Unger JB, et al. Leading topics in Twitter discourse on JUUL and Puff Bar products: content analysis. J Med Internet Res 2021 Jul 19;23(7):e26510 [FREE Full text] [doi: 10.2196/26510] [Medline: 34279236]

52. Lorentzen DG, Nolin J. Approaching completeness: capturing a hashtagged Twitter conversation and its follow-on conversation. Soc Sci Comput Rev 2015 Sep 29;35(2):277-286. [doi: 10.1177/0894439315607018]

53. Le Q, Mikolov T. Distributed representations of sentences and documents. In: Proceedings of the 31st International Conference on Machine Learning. 2014 Presented at: 31st International Conference on Machine Learning; Jun 21-26, 2014; Beijing China p. 1188-1196 URL: http://proceedings.mlr.press/v32/le14.html

54. Mikolov T, Chen K, Corrado G, Dean J. Efficient estimation of word representations in vector space. arXiv: Computer Science, Computation and Language. 2013 Jan. URL: https://arxiv.org/abs/1301.3781 [accessed 2021-09-02]

55. Saha K, Torous J, Ernala S, Rizuto C, Stafford A, De Choudhury M. A computational study of mental health awareness campaigns on social media. Transl Behav Med 2019 Nov 25;9(6):1197-1207 [FREE Full text] [doi: 10.1093/tbm/ibz028] [Medline: 30834942]

56. Zomick J, Levitan S, Serper M. Linguistic analysis of schizophrenia in Reddit posts. In: Proceedings of the Sixth Workshop on Computational Linguistics and Clinical Psychology. 2019 Presented at: Sixth Workshop on Computational Linguistics and Clinical Psychology; Jun 2019; Minneapolis, Minnesota p. 74-83. [doi: 10.18653/v1/w19-3009]

57. Devlin J, Chang M, Lee K, Toutanova K. BERT: pre-training of deep bidirectional transformers for language understanding. Assoc Comput Linguist 2019 Jun:4171-4186. [doi: 10.18653/v1/N19-1423]

58. Rogers A, Kovaleva O, Rumshisky A. A primer in BERTology: what we know about how BERT works. Trans Assoc Comput Linguist 2020 Dec;8:842-866. [doi: 10.1162/tacl_a_00349]

59. Mozafari M, Farahbakhsh R, Crespi N. A BERT-based transfer learning approach for hate speech detection in online social media. In: Complex Networks and Their Applications VIII. Basel, Switzerland: Springer International Publishing; Nov 26, 2019.

60. Heidari M, Jones J. Using BERT to extract topic-independent sentiment features for social media bot detection. In: Proceedings of the 11th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON). 2020 Presented at: 11th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON); Oct 28-31, 2020; New York, NY, USA. [doi: 10.1109/uemcon51285.2020.9298158]

61. Mozafari M, Farahbakhsh R, Crespi N. Hate speech detection and racial bias mitigation in social media based on BERT model. PLoS One 2020 Aug 27;15(8):e0237861 [FREE Full text] [doi: 10.1371/journal.pone.0237861] [Medline: 32853205]

62. Pérez-Pérez M, Pérez-Rodríguez G, Fdez-Riverola F, Lourenço A. Using Twitter to understand the human bowel disease community: exploratory analysis of key topics. J Med Internet Res 2019 Aug 15;21(8):e12610 [FREE Full text] [doi: 10.2196/12610] [Medline: 31411142]

63. Li J, Xu Q, Shah N, Mackey TK. A machine learning approach for the detection and characterization of illicit drug dealers on instagram: model evaluation study. J Med Internet Res 2019 Jun 15;21(6):e13803 [FREE Full text] [doi: 10.2196/13803] [Medline: 31199298]

64. Pereira-Kohatsu JC, Quijano-Sánchez L, Liberatore F, Camacho-Collados M. Detecting and monitoring hate speech in Twitter. Sensors (Basel) 2019 Oct 26;19(21):4654 [FREE Full text] [doi: 10.3390/s19214654] [Medline: 31717760]

65. Chakraborty S, Tomsett R, Raghavendra R, Harborne D, Alzantot M, Cerutti F, et al. Interpretability of deep learning models: a survey of results. In: Proceedings of the IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computed, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/CBDCom/IOP/SCI). 2017 Presented at: IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computed, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/CBDCom/IOP/SCI); Aug 4-8, 2017; San Francisco, CA, USA. [doi: 10.1109/uic-atc.2017.8397411]

66. Singh T, Roberts K, Cohen T, Cobb N, Wang J, Fujimoto K, et al. Social Media as a Research Tool (SMaaRT) for risky behavior analytics: methodological review. JMIR Public Health Surveill 2020 Nov 30;6(4):e21660 [FREE Full text] [doi: 10.2196/21660] [Medline: 33252345]

67. Wang X, Kou L, Sugumaran V, Luo X, Zhang H. Emotion correlation mining through deep learning models on natural language text. IEEE Trans Cybern 2020 May 12:2987064 (forthcoming). [doi: 10.1109/TCYB.2020.2987064] [Medline: 32413938]

68. Klein AZ, Sarker A, Weissenbacher D, Gonzalez-Hernandez G. Towards scaling Twitter for digital epidemiology of birth defects. NPJ Digit Med 2019 Oct 1;2:96 [FREE Full text] [doi: 10.1038/s41746-019-0170-5] [Medline: 31583284]

69. Klein AZ, Magge A, O'Connor K, Amaro JF, Weissenbacher D, Hernandez GG. Toward using Twitter for tracking COVID-19: a natural language processing pipeline and exploratory data set. J Med Internet Res 2021 Jan 22;23(1):e25314 [FREE Full text] [doi: 10.2196/25314] [Medline: 33449904]

70. Edo-Osagie O, Smith G, Lake I, Edeghere O, De La Iglesia B. Twitter mining using semi-supervised classification for relevance filtering in syndromic surveillance. PLoS One 2019 Jul 18;14(7):e0210689 [FREE Full text] [doi: 10.1371/journal.pone.0210689] [Medline: 31318885]

71. Ru B, Li D, Hu Y, Yao L. Serendipity-a machine-learning application for mining serendipitous drug usage from social media. IEEE Trans Nanobioscience 2019 Jul;18(3):324-334 [FREE Full text] [doi: 10.1109/TNB.2019.2909094] [Medline: 30951476]

72. Nikfarjam A, Sarker A, O'Connor K, Ginn R, Gonzalez G. Pharmacovigilance from social media: mining adverse drug reaction mentions using sequence labeling with word embedding cluster features. J Am Med Inform Assoc 2015 May;22(3):671-681 [FREE Full text] [doi: 10.1093/jamia/ocu041] [Medline: 25755127]

73. Strubell E, Ganesh A, McCallum A. Energy and policy considerations for modern deep learning research. AAAI 2020 Apr 03;34(09):13693-13696. [doi: 10.1609/aaai.v34i09.7123]

74.  Metropolitan and micropolitan statistical areas population totals and components of change: 2010-2019. United States Census Bureau. URL: https://www.census.gov/data/tables/time-series/demo/popest/2010s-total-metro-and-micro-statistical-areas.html [accessed 2021-08-24]

75.  Population, population change, and estimated components of population change: April 1, 2010 to July 1, 2019 (NST-EST2019-alldata). United States Census Bureau. URL: https://www.census.gov/data/tables/time-series/demo/popest/2010s-national-total.html [accessed 2021-08-24]

76.  Zhan Y, Zhang Z, Okamoto JM, Zeng DD, Leischow SJ. Underage JUUL use patterns: content analysis of Reddit messages. J Med Internet Res 2019 Sep 09;21(9):e13038 [FREE Full text] [doi: 10.2196/13038] [Medline: 31502542]

77.  Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: an imperative style, high-performance deep learning library. arXiv: Computer Science, Machine Learning. URL: https://arxiv.org/abs/1912.01703 [accessed 2021-08-24]

78.  Pennington J, Socher R, Manning C. GloVe: global vectors for word representation. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP). 2014 Presented at: Conference on Empirical Methods in Natural Language Processing (EMNLP); Oct 2014; Doha, Qatar p. 1532-1543 URL: http://www.aclweb.org/anthology/D14-1162 [doi: 10.3115/v1/d14-1162]

79.  Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Stat Soc Series B Stat Methodol 1995;57(1):289-300. [doi: 10.1111/j.2517-6161.1995.tb02031.x]

80.  Hopkins ZH, Moreno C, Secrest AM. Influence of social media on cosmetic procedure interest. J Clin Aesthet Dermatol 2020 Jan;13(1):28-31 [FREE Full text] [Medline: 32082468]

81.  Merchant RM, Asch DA, Crutchley P, Ungar LH, Guntuku SC, Eichstaedt JC, et al. Evaluating the predictability of medical conditions from social media posts. PLoS One 2019 Jun 17;14(6):e0215476. [doi: 10.1371/journal.pone.0215476] [Medline: 31206534]

82.  Eichstaedt JC, Smith RJ, Merchant RM, Ungar LH, Crutchley P, Preoţiuc-Pietro D, et al. Facebook language predicts depression in medical records. Proc Natl Acad Sci U S A 2018 Oct 30;115(44):11203-11208 [FREE Full text] [doi: 10.1073/pnas.1802331115] [Medline: 30322910]

83.  Griffis H, Asch D, Schwartz H, Ungar L, Buttenheim A, Barg F, et al. Using social media to track geographic variability in language about diabetes: analysis of diabetes-related tweets across the United States. JMIR Diabetes 2020 Jan 26;5(1):e14431 [FREE Full text] [doi: 10.2196/14431] [Medline: 32044757]

84.  Smith RJ, Crutchley P, Schwartz HA, Ungar L, Shofer F, Padrez KA, et al. Variations in Facebook posting patterns across validated patient health conditions: a prospective cohort study. J Med Internet Res 2017 Jan 06;19(1):e7 [FREE Full text] [doi: 10.2196/jmir.6486] [Medline: 28062392]

85.  Guntuku SC, Schneider R, Pelullo A, Young J, Wong V, Ungar L, et al. Studying expressions of loneliness in individuals using twitter: an observational study. BMJ Open 2019 Nov 04;9(11):e030355 [FREE Full text] [doi: 10.1136/bmjopen-2019-030355] [Medline: 31685502]

86.  Cecinati F, Matthews T, Natarajan S, McCullen N, Coley D. Mining social media to identify heat waves. Int J Environ Res Public Health 2019 Mar 02;16(5):762 [FREE Full text] [doi: 10.3390/ijerph16050762] [Medline: 30832387]

87.  Birnbaum ML, Ernala SK, Rizvi AF, De Choudhury M, Kane JM. A collaborative approach to identifying social media markers of schizophrenia by employing machine learning and clinical appraisals. J Med Internet Res 2017 Aug 14;19(8):e289 [FREE Full text] [doi: 10.2196/jmir.7956] [Medline: 28807891]

88.  Spearman C. The proof and measurement of association between two things. Am J Psychol 1904;15(1):72-101. [doi: 10.2307/1412159]

89.  Bergman BG, Wu W, Marsch LA, Crosier BS, DeLise TC, Hassanpour S. Associations between substance use and Instagram participation to inform social network-based screening models: multimodal cross-sectional study. J Med Internet Res 2020 Sep 16;22(9):e21916 [FREE Full text] [doi: 10.2196/21916] [Medline: 32936081]

90.  United States Census Bureau. Educational Attainment, Income in the Past 12 Months (In 2018 inflation-adjusted dollars), ACS Demographic and Housing Estimates. American Community Survey. 2018. URL: https://www.census.gov/acs/www/data/data-tables-and-tools/subject-tables/

91.  Behavioral risk factor surveillance system. Centers for Disease Control and Prevention. URL: https://www.cdc.gov/brfss/index.html [accessed 2021-08-24]

92.  2019 County Health Rankings Key Findings Report. University of Wisconsin Population Health Institute. URL: https://www.countyhealthrankings.org/reports/2019-county-health-rankings-key-findings-report [accessed 2021-08-24]

93.  County business patterns: 2017. United States Census Bureau. URL: https://www.census.gov/data/datasets/2017/econ/cbp/2017-cbp.html [accessed 2021-08-24]

94.  Apparent per capita alcohol consumption: national, state, and regional trends 1977-2018. Open ICPSR. URL: https://www.openicpsr.org/openicpsr/project/105583/version/V5/view;jsessionid=D71A8F94F8F63A42806C3C45BD042D9D [accessed 2021-08-24]

95.  Kavanagh A, Kelly M, Krnjacki L, Thornton L, Jolley D, Subramanian S, et al. Access to alcohol outlets and harmful alcohol consumption: a multi-level study in Melbourne, Australia. Addiction 2011 Oct;106(10):1772-1779. [doi: 10.1111/j.1360-0443.2011.03510.x] [Medline: 21615583]

96. Popova S, Giesbrecht N, Bekmuradov D, Patra J. Hours and days of sale and density of alcohol outlets: impacts on alcohol consumption and damage: a systematic review. Alcohol Alcohol 2009;44(5):500-516. [doi: 10.1093/alcalc/agp054] [Medline: 19734159]

97. Treno AJ, Johnson FW, Remer LG, Gruenewald PJ. The impact of outlet densities on alcohol-related crashes: a spatial panel approach. Accid Anal Prev 2007 Sep;39(5):894-901. [doi: 10.1016/j.aap.2006.12.011] [Medline: 17275773]

98. Belinkov Y, Gehrmann S, Pavlick E. Interpretability and analysis in neural NLP. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts.: Association for Computational Linguistics; 2020 Presented at: 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts; Jul, 2020; Online p. 1-5. [doi: 10.18653/v1/2020.acl-tutorials.1]

99. Barefoot JC, Grønbaek M, Feaganes JR, McPherson RS, Williams RB, Siegler IC. Alcoholic beverage preference, diet, and health habits in the UNC Alumni Heart Study. Am J Clin Nutr 2002 Aug;76(2):466-472. [doi: 10.1093/ajcn/76.2.466] [Medline: 12145024]

100. Sluik D, Bezemer R, Sierksma A, Feskens E. Alcoholic beverage preference and dietary habits: a systematic literature review. Crit Rev Food Sci Nutr 2016 Oct 25;56(14):2370-2382. [doi: 10.1080/10408398.2013.841118] [Medline: 25674684]

101. Dey M, Gmel G, Studer J, Dermota P, Mohler-Kuo M. Beverage preferences and associated drinking patterns, consequences and other substance use behaviours. Eur J Public Health 2014 Jun;24(3):496-501. [doi: 10.1093/eurpub/ckt109] [Medline: 23940073]

102. Alcácera MA, Marques-Lopes I, Fajó-Pascual M, Foncillas JP, Carmona-Torre F, Martínez-González MA. Alcoholic beverage preference and dietary pattern in Spanish university graduates: the SUN cohort study. Eur J Clin Nutr 2008 Oct;62(10):1178-1186. [doi: 10.1038/sj.ejcn.1602833] [Medline: 17609695]

103. Siegel MB, Naimi TS, Cremeens JL, Nelson DE. Alcoholic beverage preferences and associated drinking patterns and risk behaviors among high school youth. Am J Prev Med 2011 Apr;40(4):419-426. [doi: 10.1016/j.amepre.2010.12.011] [Medline: 21406275]

104. McCann SE, Sempos C, Freudenheim JL, Muti P, Russell M, Nochajski TH, et al. Alcoholic beverage preference and characteristics of drinkers and nondrinkers in western New York (United States). Nutr Metab Cardiovasc Dis 2003 Feb;13(1):2-11. [doi: 10.1016/s0939-4753(03)80162-x] [Medline: 12772432]

105. Lachmar EM, Wittenborn AK, Bogen KW, McCauley HL. #MyDepressionLooksLike: examining public discourse about depression on Twitter. JMIR Ment Health 2017 Oct 18;4(4):e43 [FREE Full text] [doi: 10.2196/mental.8141] [Medline: 29046270]

106. Alberga A, Withnell S, von Ranson KM. Fitspiration and thinspiration: a comparison across three social networking sites. J Eat Disord 2018 Nov 26;6:39 [FREE Full text] [doi: 10.1186/s40337-018-0227-x] [Medline: 30534376]

107. Shepherd A, Sanders C, Doyle M, Shaw J. Using social media for support and feedback by mental health service users: thematic analysis of a twitter conversation. BMC Psychiatry 2015 Feb 19;15:29 [FREE Full text] [doi: 10.1186/s12888-015-0408-y] [Medline: 25881089]

108. van Swol LM, Chang C, Kerr B, Moreno M. Linguistic predictors of problematic drinking in alcohol-related Facebook posts. J Health Commun 2020 Mar 03;25(3):214-222 [FREE Full text] [doi: 10.1080/10810730.2020.1731632] [Medline: 32096449]

109. Park A, Conway M, Chen AT. Examining thematic similarity, difference, and membership in three online mental health communities from reddit: a text mining and visualization approach. Comput Hum Behav 2018 Jan;78:98-112 [FREE Full text] [doi: 10.1016/j.chb.2017.09.001] [Medline: 29456286]

110. Kavuluru R, Williams AG, Ramos-Morales M, Haye L, Holaday T, Cerel J. Classification of helpful comments on online suicide watch forums. ACM BCB 2016 Oct;2016:32-40 [FREE Full text] [doi: 10.1145/2975167.2975170] [Medline: 28736770]

## Abbreviations

**API:** application programming interface
**BERT:** Bidirectional Encoder Representations From Transformers
**MMSA:** metropolitan-micropolitan statistical area
**NLP:** natural language processing

XSL•FO
**RenderX**