

BMJ Open Simplified data science approach to extract social and behavioural determinants: a retrospective chart review

Andrew Teng , Adam Wilcox

To cite: Teng A, Wilcox A. Simplified data science approach to extract social and behavioural determinants: a retrospective chart review. *BMJ Open* 2022;**12**:e048397. doi:10.1136/bmjopen-2020-048397

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2020-048397>).

Received 19 January 2021
Accepted 07 November 2021



© Author(s) (or their employer(s)) 2022. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

Biomedical Informatics and Medical Education, University of Washington, Seattle, Washington, USA

Correspondence to

Andrew Teng; akteng@uw.edu

ABSTRACT

Objectives We aim to extract a subset of social factors from clinical notes using common text classification methods.

Design Retrospective chart review.

Setting We collaborated with a local level I trauma hospital located in an underserved area that has a housing unstable patient population of about 6.5% and extracted text notes related to various social determinants for acute care patients.

Participants Notes were retrospectively extracted from 43 798 acute care patients.

Methods We solely use open source Python packages to test simple text classification methods that can potentially be easily generalisable and implemented. We extracted social history text from various sources, such as admission and emergency department notes, over a 5-year timeframe and performed manual chart reviews to ensure data quality. We manually labelled the sentiment of the notes, treating each text entry independently. Four different models with two different feature selection methods (bag of words and bigrams) were used to classify and predict housing stability, tobacco use and alcohol use status for the extracted clinical text.

Results From our analysis, we found overall positive results and metrics in applying open-source classification techniques; the accuracy scores were 91.2%, 84.7%, 82.8% for housing stability, tobacco use and alcohol use, respectively. There were many limitations in our analysis including social factors not present due to patient condition, multiple copy-forward entries and shorthand. Additionally, it was difficult to translate usage degrees for tobacco and alcohol use. However, when compared with structured data sources, our classification approach on unstructured notes yielded more results for housing and alcohol use; tobacco use proved less fruitful for unstructured notes.

INTRODUCTION

Most data can be generally categorised as structured or unstructured, where structured data can consist of items such as vital signs and lab results and unstructured data can consist of items such as text notes or images.¹ Although structured data can generally be easier to extract and analyse, unstructured

Strengths and limitations of this study

- From our analysis, we can first see that text classifiers are promising when applied to extracted clinical notes for housing stability, tobacco use and alcohol use status.
- Additionally, we found that structured data sources, such as diagnosis codes and intake surveys, vary and may not be the most holistic approach to understanding housing stability, tobacco use and alcohol use.
- Our simplified approach has shown that open source simple text classifiers can be used to predict text sentiment for social and behavioural determinants and can supplement current structured sources to provide a more complete social history for patients.
- However, even with a few limitations with our approach, we believe that this workflow can help inform clinicians and provide an easily implementable snapshot on patient social history.

data can potentially provide an array of information not present or easily identifiable in structured data. As healthcare institutions expand data collection to include non-clinical features, more unstructured data surrounding behavioural health and social determinants of health (SDoH) information are starting to become more readily available. Furthermore, there has been a growing interest around Medicaid patients, as SDoH can drive up to 80% of health outcomes, especially within this patient demographic.² Therefore, SDoH and REAL (Race, Ethnicity and Language) data are now being used for secondary analysis as recent research has indicated that there is a correlation between SDoH and health outcomes and the increasing need to research health disparities across populations.³

SDoH and REAL can include housing stability, access jobs and healthcare services, education level, language and socioeconomic conditions.⁴ These indicators are descriptors



of different societies and are useful as predictors of health outcomes and the uptake of health interventions.⁵ Because they can potentially be powerful indicators of health, many institutions are now starting to analyse and intake SDoH and REAL information, whether through text notes or standardised coding, such as International Classification of Diseases (ICD).⁶ Additionally, SDoH can provide health teams with a greater understanding of a patient condition holistically.⁷ However, there are challenges with SDoH intake as there is no standardised SDoH screening tool in the Electronic Health Record (EHR) across institutions⁸; additionally, coding schemes like ICD can prove to be unreliable in secondary analysis as coding can oversimplify symptoms and diagnoses leading to coding uncertainties and the fact that coding errors may be present from unintentional mistakes or even upcoding.^{9,10} Additionally, certain SDoH data may be more complete than others due to reimbursement incentives or other priorities.¹¹ Past research has shown that hospital readmissions are highly influenced by patient health status and SDoH and suggest that clinical staff and researchers should consider SDoH when assessing readmission risk.¹²

The 2018–2019 King Country Community Health Needs Assessment (CHNA) reported the results from a health need assessment survey given to residents to identify regional perceived healthcare issues. It was determined that housing affordability and housing stability were major challenges dominating overall health.¹³ Mental health was also highlighted as a challenge for healthcare providers; mental illness can be caused by depression, schizophrenia and alcohol and substance-related disorders.¹³ The CHNA reported that adults in the lowest income tier were about 15 times more likely to experience severe psychological distress compared with their high-income counterparts. Additionally, it is noted that part of the region had continued challenges with adult smoking rates.¹³ Locally, it is estimated that there are at least 22 000 homeless individuals in King Country and more than 12 000 people in the Seattle region, a 4% increase over the previous year.¹⁴ Housing instability is associated with various health inequalities, such as shorter life expectancy, higher morbidity and increased usage of acute hospital services, ‘as the social determinants of homelessness and health inequities are often intertwined, and long-term homelessness further exacerbates poor health’.¹⁵ It is, therefore, important to treat housing stability and other SDoH as a combined health issue to aid in improving health outcomes in clinical settings. Although some research have shown that patients who experience housing instability are more likely to die following admission for severe sepsis than those with insurance,¹⁶ other research indicates that the effects of health inequalities are still unclear and need further investigation.¹⁷ Additionally, various behavioural habits, including tobacco and alcohol use, although may not directly be considered a SDoH, can impact health decisions and outcomes. For example, one study found that

participants who drank alcohol and reported tobacco use consumed more foods higher in fat and sugar, low in vitamins and minerals as well as foods, considered by them to be less healthy and prepared in a less healthy way.¹⁸

Within our region, it has been noted in recent years that the smoking rate is around 13%; however, among Black/African-Americans or individuals with multiple races, is double the rate among white adults and four times higher than Asian adults. Additionally, it was reported that, when compared with high-income households, low-income households were three times more likely to be smokers.^{13,19} Drug and alcohol use also shared similar metrics; within the region, ‘drug and alcohol-caused deaths was 22% higher among Blacks and four times greater among American Indian/Alaskan Native than among non-Hispanic Whites’ and alcohol use represented 4.97/100 000 deaths locally in 2015.^{20,21} Therefore, it may be important to look at social determinants and health behaviours, together known as social and behavioural determinants of health (SBDH) to better understand the patient population.¹⁸

Recent technological advances in machine learning and artificial intelligence have shown great potential in providing a pathway for informaticians and clinicians to better understand unstructured data.

Within the clinical setting, there have been numerous approaches in adopting natural language processing (NLP) to aid with processing unstructured clinical text notes. Common uses of NLP include extracting diagnoses and chief reports as well as grouping of information for quality improvement. There are various NLP methods that can be used in the clinical setting, such as automatic tagging of conditions or variables of interest, sentiment classification or even text extraction. Various open source NLP and ontological tools, such as Automated Retrieval Console, Apache clinical Text Analysis and Knowledge Extraction System (Apache cTAKES), MetaMap and HITEX, Unified Medical Language System (UMLS) Metathesaurus and BioPortal have been used to aid with text extraction or classification.^{22–24} On the other hand, less complex classification methods have been used as well to identify specific groups of patients, risk assessment or aid in validating structured annotation.^{25–27} A recent scoping review found that although practitioners collect a variety of SBDH data at point of care through EHR, the overall use of automated technology is limited to date.²⁸

With the idea of implementing an easily generalisable approach to classify selected social factors, we extracted both unstructured and structured data sources related to SBDH from a local hospital to identify and generate a framework to automatically extract and classify SBDH from text notes. We focused on housing stability status, tobacco use and alcohol use. These three social factors were chosen due to their direct impact on health outcomes and the local public health impact^{14–18} and presence in the EHR. To tackle challenges associated with SBDH extraction from unstructured text notes, we aimed to create a generalisable framework using low barrier open-source tools that are commonly used in

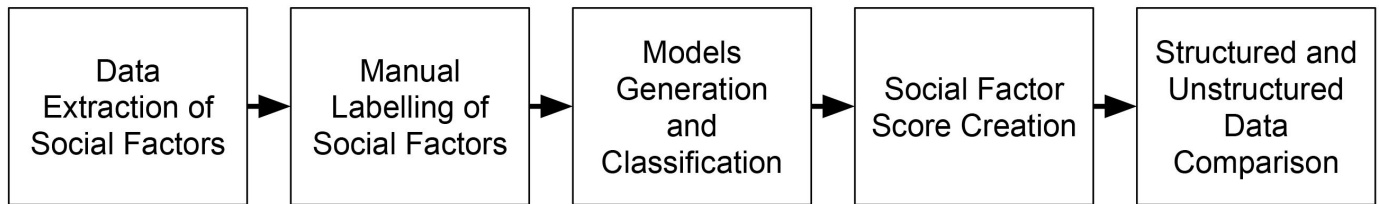


Figure 1 High-level overview of the workflow process.

the data science field. Because notes and stylistic choices can be institution and location specific, we sought not to create a model that is generalisable but rather a simplified method that could be potentially easily implemented using common off the shelf NLP and data science tools.

METHODS

Study design and overview

A high-level overview of our workflow is seen in [figure 1](#). We retrospectively extracted patient data from the acute care setting at a level I trauma centre and academic teaching hospital with the aim to create a general and easily applicable workflow to extract and classify SBDH factors from clinical notes. We applied a two-pronged approach and collected unstructured data from a subset of patients over a 1-year timespan (group A) to create and test the text classification model and also collected structured and unstructured data from a subset of patients over a 5-year timespan (group B) to apply the best model created from group A and compare results between the two data types. We performed automatic classification and scoring of patients via various NLP classification methods on three social factors: (1) housing stability, (2)

tobacco use and (3) alcohol use. Our general workflow for housing stability, a similar approach was also used for tobacco and alcohol use, is seen in [figure 2](#).

Study population

Data were not only extracted from Harborview Medical Center, a 413-bed academic hospital that has a patient population consisting mostly from Washington, but also from a five-state area.²⁹ In 2014, there were 17 121 inpatient admissions, where 19% of the patients belong to a racial or ethnic minority and 37% of patients were enrolled in Medicaid.^{29 30} Additionally, in 2015, the non-US born population was estimated to be around 21% in Seattle highlighting the potential diversity that could be found with this patient population.³⁰

Data sources, extraction and validation

We extracted both structured and unstructured data sources related to housing stability, tobacco use and alcohol use using Structured Query Language (SQL) queries called directly from an integrated python-based Jupyter Notebook:

1. Structured data sources include billing and diagnostic/ICD 9 and 10 codes, questionnaire or Epic SmartForm

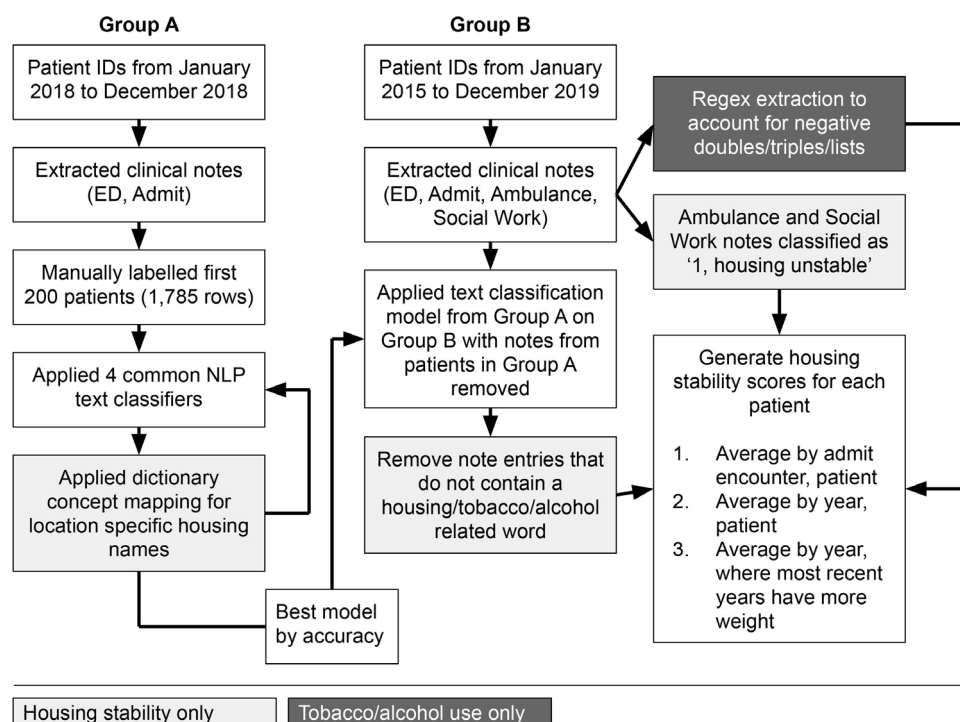


Figure 2 Text extraction, classification and scoring workflow. ED, emergency department.

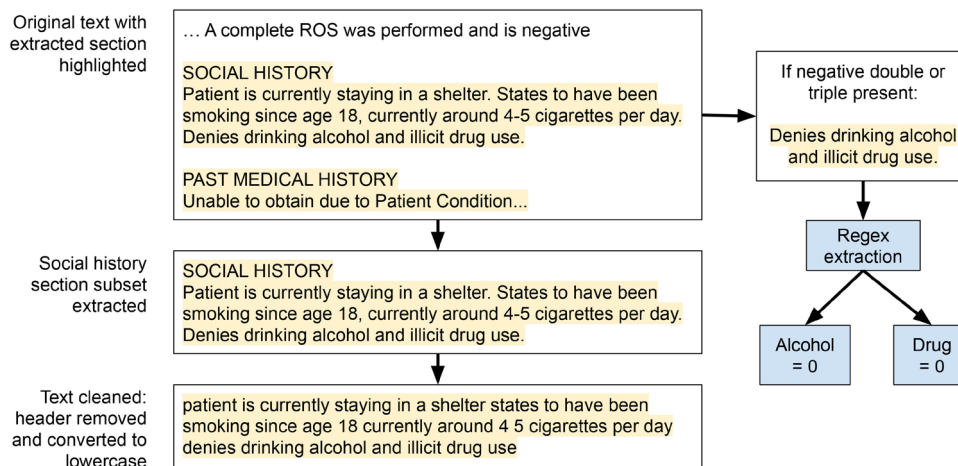


Figure 3 Text extraction and cleaning process. Additional steps were performed for notes when classifying text related to tobacco and alcohol use to extract negative sentiment doubles or triples. ROS, Review of Systems.

responses, address fields (location), problem list (ICD 9), patient encounters, clinical events (actual encounters of care) and discharge/disposition location.

2. Unstructured data sources consisted of text notes from the emergency department (ED), admission (admit) notes, social work and ambulance notes.

Discharge notes were not explored as they were not recorded in the same subdivided format as the admit and ED notes, making selective text extraction of SBDH difficult. From our initial list of patient identifiers over a 1-year timespan from group A, we performed manual EHR validation of a random subset of 50 patients to validate the completeness of the clinical notes and confirm the location of social history and social factors in clinical notes. Extensive research and conversations with an internal data analyst confirmed the location of these topics (housing, tobacco use and alcohol use) within structured data sources.

Data cleaning

After confirmation, clinical notes were extracted for both groups A and B. The notes were cleaned (eg, symbols removed, converted to lowercase) prior to classification and analysis in the Python Jupyter notebook via Natural Language Toolkit (NLTK). Our general text extraction and cleaning workflow can be seen in figure 3. However, housing stability notes and tobacco or alcohol use notes were stylistically and grammatically different, and both sets needed distinct additional cleaning steps. Housing stability notes that contained the phrase ‘not homeless’ were converted via regex to say ‘housed’ instead. Additionally, for housing stability, a concept dictionary was also created to substitute local facility names with more general concept (eg, ‘Union Gospel Mission’ was converted to ‘shelter’). This was done to explore how the algorithms handle formal nouns.

For text notes in group B, we performed an additional concept extraction step. Tobacco use and alcohol use notes often contained incomplete (lacking the subject, predicate, object format) triples or doubles (eg, ‘Denies

smoking, drinking, drugs’). Due to their incomplete sentence structures, common NLP tools to parse, extract and classify triples, such as Stanford CoreNLP, were not suitable as these tools rely on having all three parts of the triple present. These notes related to tobacco and alcohol use, therefore, underwent an additional step that performed a separate relation extraction that would first identify a negative sentiment word (eg, denies), then individually extract the following SBDH-related objects in the list by commas or conjunctions (eg, and, or), and then label, or reclassify if necessary, the negative sentiment to all components of the list. Our process is seen in the left side of figure 3. If the regex extraction of negative lists resulted in a different result from the text classification prediction, the regex extraction would overwrite the end result prior to scoring. Once these steps were performed, the data were considered clean and suitable for classification.

Model building

Cleaned text from group A were used to generate and test the classification models. These notes were split in 70/30 validation and testing sets. We applied four different common NLP text classification models to the testing sets (via SciKit Learn): multinomial naïve Bayes, support vector machine, logistic regression and random forest. Default parameters and a bag-of-words approach were used. The best-performing model by accuracy was then chosen and applied to the larger corpus, Group B, with notes from patients in Group A removed, to avoid overfitting and classification bias. This process was performed for housing, tobacco use and alcohol use.

Scoring generation

In order to create a simple method of identifying patients who are experiencing social instability, we created a scoring metric based on the classified notes. After applying the optimum model by accuracy to the entire corpus of extracted text notes, housing stability, tobacco use and alcohol use scores were generated. Patient

identifiers were mapped by patient location and those who were not in the acute care setting during this time-frame were removed. Three different scoring approaches were used to describe these social factors: (1) predictions were averaged by patient encounter, then averaged by patient identifier, (2) predictions were averaged by year, then by patient identifier and (3) predictions were averaged by year, where each year then had a weight where the most recent year had the highest weight and the furthest year had the lowest weight (eg, predictions from 2019 were weighted by a factor of 5 and predictions from 2015 were weighted by a factor of 1). This scoring generation process was then repeated on our structured data for all three social factors and the results were compared and analysed. Structured data were also extracted for our list of patients in group B.

Patient and public involvement

No patients were involved. The retrospective exploration is a part of a larger study and was approved by the University of Washington Institutional Review Board #STUDY00006723. Patient data elements, including encounter identifiers, race, age and notes with SBDH, were extracted directly from the data warehouse and stored on encrypted computers and were not distributed or shared outside of the secured and closed environment. No patient identifiers or names were stored in this analysis.

RESULTS

Characteristics of study subjects

Clinical notes (ED, admit, social work and ambulance) between 2015 and 2019 were extracted and included, forming group B. Notes from the first 200 patients were included in group A and notes from 1 47 457 patients were included in group B. During the same time frame, 61 767 patients were in acute care. After extraction and model prediction, the patient notes were cross-referenced with inpatient location and only notes from those who were in acute care were retained, for a total of 43 798 patients from 2015 to 2019. The patient demographics of this final subset were 63% (n=27 575) men, 37% (n=16 223) women, 88.2% (n=38 634) not Hispanic or Latino and 10.5% (n=4 609) Hispanic or Latino and 1.3% (n=555) unknown or not answered. Further descriptive statistics are found in [table 1](#).

Data attributes

[Table 2](#) illustrates the amount of data for each corresponding extraction level, specifically for housing status. We first started with extracting text from the ED and admit notes, forming group A, which consisted of 50 000 rows or text entries and covered 3 200 unique patients, over a 1-year time frame. From there, we manually labelled housing stability concepts in a binary fashion, where 0 would indicate housing stability and 1 would indicate any level of housing instability, regardless of severity. As

Table 1 Population demographics

Race (n=43 798)	n (%)
White or Caucasian	31 575 (72.1)
Black or African American	4 812 (11.0)
Asian	3 174 (7.2)
American Indian or Alaska Native	1 165 (2.7)
Native Hawaiian or other Pacific Islander	524 (1.2)
Multiple races	3 (0)
Unavailable, unknown or missing	2 545 (5.8)
Age range (n=43 798)	n (%)
0–18	1 856 (4.2)
19–44	12 437 (28.4)
45–64	14 863 (33.9)
65–84	11 902 (27.2)
85 and over	2 740 (6.3)

manual labelling can be a labour-intensive process, only the first 6 000 text rows were labelled, covering 218 unique patients. However, within these first 6 000 rows, numerous notes did not contain text that alluded to housing status or were empty due to patient condition. Therefore, only 1 785 out of the 6 000 rows were labelled, covering 200 unique patients, where 995 (55.7%) were labelled as housing stable and 790 (44.3%) were labelled as housing unstable. We also found that 5.7% of the entries within this subset were duplicates or copy-forward entries. The same workflow was performed for labelling tobacco and alcohol use. However, only 1 108 rows were labelled for tobacco use and 1 220 rows for alcohol use, where in both cases, 0 indicated no use, 1 indicated rare/previous/occasional use and two indicated current use, regardless of degree. Tobacco use resulted in 446 (40.3%) labels for no use, 129 (11.6%) labels for rare/previous/occasional use and 533 (48.1%) labels for current use. Similarly, alcohol use resulted in 595 (48.8%) labels for no use, 185 (15.2%) labels for rare/previous/occasional use and 440 (36%) labels for current use.

Model performance

Four different common text classifiers, mentioned in the Methods section, were applied to the manually labelled group A data. The statistical metrics, including accuracy, precision and recall, are seen in [tables 3 and 4](#). The accuracies between the classifiers and each classification technique for housing stability were overall fairly high ranging from 84.36% to 92.18%. The accuracies for tobacco and alcohol use were lower, ranging from 70.87% to 84.68% for tobacco use and 69.95% to 82.79% for alcohol use. Additionally, for each top performing model, the most influential words for text classification, for each social factor, are seen in [table 5](#). The best-performing

Table 2 Extracted data amounts for housing status

Level of extraction	Rows (n)	Unique patients (n)	Unique encounters (n)	Social history entries (n/unique)
ED and admit notes	49955	3233	15664	21 876/21334
Housing, tobacco, alcohol information	6000	218	1995	2408/2211
Remove nulls/missing data	Housing: 1785 Tobacco: 1108 Alcohol: 1220	Housing: 200 Tobacco: 179 Alcohol: 181	1361	1785/1684

ED, emergency department.

classification models were selected for each social factor and were used to apply the model to our entire corpus in group B.

Scoring results and comparison

After classifying text for housing stability, tobacco use and alcohol use for patients in group B, we applied a scoring metric scheme, described in the Methods section. We generated three different scores that were calculated and weighted differently based on time. Our final score weighs more recent note entries and their resulting classification score higher than notes from previous years as social factors and their influence can change over time. Using the same process, we extracted and scored housing stability, tobacco use and alcohol use with structured data sources and compared the results with the unstructured process.

Housing stability

Using notes, we classified 839 patients as housing unstable, a score above 0.5, and 21 370 patients as housing stable, a score of 0.5 and below. In total, we classified 22 209 patients with this text classification workflow, which covered 50.71% of the acute care patients within the same timeframe. When compared with structured data sources, only 791 (1.81%) additional patients were found.

Table 3 Accuracies among text classifiers

	n=1	n=1-2
Multinomial naïve Bayes	Housing: 91.62% Tobacco: 70.87% Alcohol: 70.77%	Housing: 91.43% Tobacco: 77.18% Alcohol: 69.95%
Support vector machine	Housing: 92.18% Tobacco: 81.08% Alcohol: 76.50%	Housing: 91.99% Tobacco: 82.88% Alcohol: 81.97%
Logistic regression	Housing: 84.36% Tobacco: 75.38% Alcohol: 77.60%	Housing: 90.13% Tobacco: 84.68% Alcohol: 82.79%
Random forest	Housing: 90.50% Tobacco: 76.28% Alcohol: 71.31%	Housing: 91.25% Tobacco: 78.98% Alcohol: 75.68%

Bold values indicates highest performance for each SDBH.

Tobacco use

We classified 4911 patients as currently using tobacco, regardless of amount or degree (1.5–2) using text notes. We classified 1480 patients as having rare/occasional/past use of tobacco (0.5–1.5), and 7139 patients as not using tobacco (0–0.5). In total, we classified 13 530 patients with this text classification workflow, which covered 30.9% of the acute care patients within the same timeframe. When compared with structured data sources, 179 351 (40.9%) additional patients were captured.

Alcohol use

We classified 2738 patients as currently using alcohol, regardless of amount or degree (1.5–2) using text notes. We classified 4050 patients as having rare/occasional/past use of alcohol (0.5–1.5), and 13 885 patients as not drinking alcohol (0–0.5). In total, we classified 20 673 patients with this text classification workflow, which covered 37% of the acute care patients within the same timeframe. When compared with structured data sources, no additional patients were found.

DISCUSSION

Our approach to a simple text classification method for various SDoH has shown positive results. The selected classification models were chosen as they were the most commonly used classification models when researching text classification techniques. Furthermore, these models were robust enough to curtail the need for more complex machine learning-based text classification methods, which may be harder to interpret in the clinical space as the weights and decisions can be confiscated due to the black box nature of these more complex classification methods. In general, linear models are fast to train, can work well with sparse data and offer interpretability.³¹ Additionally, recent research has also suggested that more complex machine learning approaches may not yield statistically significant improvements in predictive power to justify the time and effort necessary to implement and test these more complex methods. Although promising, more advanced methods of NLP, such as convoluted neural networks, may not provide a significant tradeoff in improvement or accuracy versus transparent understanding of rule-based approaches. In fact,

Table 4 Best-performing classifier detailed metrics

	Classifier	Accuracy	Recall	Precision	F1
Housing status*	Support vector machine (n=1)	0.92	0.93/0.91 (0/1)	0.94/0.90	0.93/0.91
Tobacco use†	Logistic regression (n=1–2)	0.85	0.82/0.95/0.86 (0,1,2)	0.96/0.43/0.87 (0,1,2)	0.88/0.60/0.87 (0,1,2)
Alcohol use†	Logistic regression (n=1–2)	0.83	0.86/0.73/0.81 (0,1,2)	0.93/0.44/0.88 (0,1,2)	0.89/0.55/0.84 (0,1,2)

*0: no use, 1: current use.

†0: no use, 1: rare/occasional/history, 2: current use.

Yao *et al* found that the F1 scores for Convolutional neural network (CNN) via TensorFlow did not improve significantly for interested features when compared with logistic regression and support vector machine implementations.³² Finally, generalisable methods to create institution-specific models can be better for the health-care system as a whole as each institution records clinical information with variances.

Although SBDH information and other social factors can be indicative of overall health, collection of SBDH heavily relies on clinical staff to screen and document SBDH. Furthermore, it also assumes that patients will respond accurately and truthfully. Various financial incentives from the federal level have propelled collection of social factors, such as tobacco use and tobacco cessation. However, other social factors, which can be equally as important, such as alcohol use, are not incentivised to be captured; rather only more severe instances are incentivised, such as alcohol dependence or alcohol addiction or disorder.^{11 33} Due to this discrepancy, we found that structured data sources were less reliable, and that text classification aided in detailing a patient more holistically.

Our text classification of unstructured data relied solely on ED, admit, social work and ambulatory notes as our parsing and extraction method could only work with notes in a certain format with the social history heading. Social factors and other social history could also be recorded in other locations but were not compatible with our approach. Furthermore, social work and

ambulatory notes used for housing status only and were only extracted if the notes contained a word or phrase related to housing instability. This approach was used as the notes were typically stored in a more unstructured format compared with the ED and admit notes; there were no section headers. The lack of section headers increased the difficulty to extract the notes and the notes would often verbiage that would interfere with the simple text classification approach that we used. Therefore, we decided to extract notes that contained words relating to housing instability. Additionally, tobacco and alcohol use notes had stylistic and grammatical challenges. These social factors were often grouped together in incomplete triples (eg, 'denies drinking, smoking, illicit drug use'). The classification algorithms often had trouble reciprocating the negative connotation to all components of the triple. Therefore, we used regex to specifically extract these triples and classify the note based on the presence of words related to tobacco or alcohol. Without this additional data cleaning or manipulation step, the negative sentiment in a list would not have been applied to all elements within the list, but rather only the first element. In our example of 'denies smoking, drinking, drugs', the negative sentiment of 'denies' would have only been applied to smoking as smoking immediately follows 'denies'. However, with our additional concept extraction step, the negative sentiment of 'denies' is now also applied to 'drinking' and 'drugs'. These results would then override the text classification algorithm, if

Table 5 Word or phrase importance ranking

Social factor (classifier)	Top 20 weighted words
Housing stability (support vector machine, n=1)	['friends' 'motel' 'stay' 'cigs' 'found' 'street' 'stays' 'streets' 'van' 'incarcerated' 'desc' 'currently' 'undomiciled' 'friend' 'respite' 'kcj' 'shelters' 'homelessness' 'shelter' 'homeless']
No tobacco use (logistic regression, n=1,2)	['use denies' 'denies' 'lives' 'tobacco drug' 'seattle denies' 'use results' 'lives seattle' 'alcohol tobacco' 'tobacco drugs' 'never smoker' 'etoh tobacco' 'drinking' 'seattle tobacco' 'denies cigarettes' 'drugs tobacco' 'denies alcohol' 'tobacco alcohol' 'denies smoking' 'denies' 'denies tobacco']
No alcohol use (logistic regression, n=1,2)	['care' 'ppd' 'tobacco' 'smoking' 'etoh tobacco' 'history cocaine' 'tobacco alcohol' 'etoh illicit' 'alcohol tobacco' 'etoh drug' 'drugs etoh' 'alcohol drug' 'use none' 'alcohol drugs' 'drug etoh' 'denies alcohol' 'lives' 'denies drug' 'denies etoh' 'denies']

there was a discrepancy. Therefore, the scoring metrics for these cases would not necessarily reflect the accuracy or performance of our scoring method.

It was interesting to find that tobacco use was recorded significantly more often in structured data sources compared with alcohol use and housing stability. However, because tobacco use is a (Centres for Medicare and Medicare Services) CMS core quality measure, it can be expected that this feature is more available in structured form as it is often directly asked to the patient on intake forms, screeners or during cessation treatment.¹¹ Furthermore, the Joint Commission created the Tobacco Performance Measure Set, which are three standardised performance measures addressing tobacco screening and cessation counselling: (1) tobacco use screening of patients 18 years and over, (2) tobacco use treatment, including counselling and medication during hospitalisation and (3) tobacco use treatment management plan at discharge. CMS began using these performance measures in 2016.³⁴ Because alcohol consumption is not a recommended CMS core quality measure for adults, the amount of data regarding alcohol use is not complete in structured form as it may not be consistently collected during intake procedures.

Past research has consistently pointed towards SBDH impacting patient health and outcomes. However, collection of SBDH can be a major limiting factor in the ability to model and integrate these data. There has not been a standardised collection process for SBDH data across the institution, whether it is recorded through notes or electronic forms. Additionally, many times, SBDH data may not be asked due to patient condition or it might not be updated regularly. Providers and healthcare institutions should strive to collect SBDH data more regularly even if the data fields are not empty as SBDH status can change. These intake procedures should be present and not optional; currently, only language preference must be completed due to translation laws in place. Additionally, educating patients to use patient portals and update information via these portals can provide more current SBDH information. However, we should note that vulnerable populations would most likely not be the primary audience to use this feature, and this is the subpopulation that arguably needs more attention.

Limitations

Our study has numerous limitations. There were two distinct areas in our workflow that required manual attention: (1) EHR review and (2) labelling of features. Manual EHR review was performed to ensure that the notes contained social history information in a consistent location prior to widespread text extraction. We initially validated this with a random set of 10 patients, but later expanded our validation to 25 patients. We felt that having consistent results with the 25 patients indicated a high level of confidence. Manual labelling of features was time-consuming and taxing. Although only one author performed the feature labelling, having multiple team

members would provide better and possibly more consistent classification.

This approach, although we aim to create a generalisable workflow, is still stunted by local customisations due to unique nuances in note-taking language. Patients can withhold information about their social challenges, making text classification harder to perform due to incorrect incoming data streams. Our approach relies on the fact that the patient has been seen within the healthcare system at some point in the past 5 years. This approach would not be applicable to those who are new to the institution or those who are not immediately identifiable. Classification levels for unstructured notes are not concrete as descriptive wording is also not concrete and can vary (eg, 'patient was a former smoker', 'patient quit last week', 'patient is an occasional smoker', etc). Structured data sources can add a more concrete sense to the classification. There were 5.7% copy-forward entries present as data collection of social factors may not always be appropriate (eg, patient is inebriated, in an altered mental state, etc). We did not incorporate outside ontologies, such as UMLS or MetaMap, as we were interested in creating a simple text classification approach that did not need to rely on outside entities. Furthermore, we believe that these ontologies would not have added a significant improvement in our approach due to the social factors (housing, alcohol, tobacco) that were investigated. Although minimised, applying NLP to clinical notes will always present limitations and risks with biased models, biased data and data privacy.³⁵

Community needs are constantly changing as the health of the community is not static. Currently, the King County CHNA has identified obesity, healthcare access, insurance status and drug use as other potential SBDH information to explore. These data types would be stored in different areas of the EHR and within different notes. It would be interesting to see if our designed workflow presented could be applicable and generalised to meet the needs of other SBDH data. Although we aimed to create a simplified framework to extract SBDH data from clinical notes, more complex methods such as convoluted neural networks and more advanced NLP part of speech tagging may be worth exploring as they may help improve accuracy and precision of the classification. As more notes become available for patients, it will also be important to keep in mind the potential bias of having more notes present from sicker patients and evaluating ways to reduce this bias.

We sourced data from solely one medical centre. Patients might have had encounters or other visit types in neighbouring hospitals and healthcare systems in the region. The lack of data sharing between institutions prevents holistic collection of SBDH data. Data completeness is vitally important to the quality and accuracy of models that are dependent on big data. Poor data quality and completeness lead to lower utilisation and the lack of data can potentially lead to mistakes in the decision-making process; additionally, since there is no single or

standardised source for SBDH data, the diversity of data and complexity of the associated data structures increase the difficulty and bottlenecks for data integration.³⁶ The lack of a standardised methodology to collect and store all SBDH data will limit the potential of this research field. Additionally, SBDH factors are constantly changing for patients as their behaviours can change depending on their circumstance. Being able to aggregate these data and create adaptable models is crucial as these features are never static. Furthermore, public health and outreach services fluctuate over time. Creating a method or using an Application programming interface (API) to update the list of community shelters and other places for homeless services would be necessary to maintain an accurate understanding of a patients' housing status.

CONCLUSION

From our analysis, we can first see that text classifiers are promising when applied to extracted clinical notes for housing stability, tobacco use and alcohol use status. Additionally, we found that structured data sources, such as diagnosis codes and intake surveys, vary and may not be the most holistic approach to understanding housing stability, tobacco use and alcohol use. Our simplified approach has shown that open source simple text classifiers can be used to predict text sentiment for social determinants and can supplement current structured sources to provide a more complete social history for patients. However, even with a few limitations with our approach, we believe that this workflow can help inform clinicians and provide an easily implementable snapshot on patient social history.

Correction notice This article has been corrected since it first published. The areas redacted in the previous version have now been added.

Acknowledgements Sally Lee, Abdelhak Abdou, Marion Granich, David Carlborn

Contributors AT performed the data extraction, tool building and analysis. AW provided guidance and verification when needed. AT is the guarantor.

Funding This work was supported by the U.S. Department of Health and Human Services, National Library of Medicine Training Grant T15LM007442.

Competing interests None declared.

Patient and public involvement Patients and/or the public were not involved in the design, or conduct, or reporting, or dissemination plans of this research.

Patient consent for publication Not applicable.

Ethics approval This study does not involve human participants.

Provenance and peer review Not commissioned; externally peer reviewed.

Data availability statement No data are available. The data used are unable to be shared due to patient privacy, confidentiality, and US healthcare laws.

Open access This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

ORCID iD

Andrew Teng <http://orcid.org/0000-0002-6565-8239>

REFERENCES

- Raghupathi W, Raghupathi V. Big data analytics in healthcare: promise and potential. *Health Inf Sci Syst* 2014;2:3.
- Hood CM, Gennuso KP, Swain GR, *et al.* County health rankings. *Am J Prev Med* 2016;50:129–35.
- Gottlieb LM, Tirozzi KJ, Manchanda R, *et al.* Moving electronic medical records upstream. *Am J Prev Med* 2015;48:215–8.
- Social Determinants of Health. Social determinants of health. Available: <https://www.healthypeople.gov/2020/topics-objectives/topic/social-determinants-of-health> [Accessed 1 Feb 2020].
- Social Determinants. Institute for health metrics and evaluation. Available: <http://www.healthdata.org/social-determinants> [Accessed 1 Feb 2020].
- Nerenz DR. Health care organizations' use of race/ethnicity data to address quality disparities. *Health Aff* 2005;24:409–16.
- Andermann A, CLEAR Collaboration. Taking action on the social determinants of health in clinical practice: a framework for health professionals. *CMAJ* 2016;188:E474–83.
- Olson DP, Oldfield BJ, Navarro SM. Standardizing social determinants of health assessments, 2019. Available: <https://www.healthaffairs.org/doi/10.1377/hblog20190311.823116/full/>
- Wockenfuss R, Frese T, Herrmann K, *et al.* Three- and four-digit ICD-10 is not a reliable classification system in primary care. *Scand J Prim Health Care* 2009;27:131–6.
- O'Malley KJ, Cook KF, Price MD, *et al.* Measuring diagnoses: ICD code accuracy. *Health Serv Res* 2005;40:1620–39.
- Eligible professional meaningful use core measures measure 9 of 13, 2014. Available: https://www.cms.gov/Regulations-and-Guidance/Legislation/EHRIncentivePrograms/downloads/9_Record_Smoking_Status.pdf
- Lax Y, Martinez M, Brown NM. Social determinants of health and hospital readmission. *Pediatrics* 2017;140:e20171427.
- King County community health needs assessment 2018/2019. Available: <https://www.kingcounty.gov/depts/health/data/community-health-indicators/~media/depts/health/data/documents/2018-2019-Joint-CHNA-Report.ashx>
- Henry M, Mahathay A, Morrill T. The 2018 annual homeless assessment report (AHAR) to Congress. The U.S. department of housing and urban development office of community planning and development, 2018. Available: <https://files.hudexchange.info/resources/documents/2018-AHAR-Part-1.pdf>
- Stafford A, Wood L. Tackling health disparities for people who are homeless? Start with social determinants. *Int J Environ Res Public Health* 2017;14:1535.
- Ahmad S, Baig S, Taneja A, *et al.* The outcomes of severe sepsis in homeless. *Chest* 2014;146:230A.
- Bambra C, Gibson M, Sowden A, *et al.* Tackling the wider social determinants of health and health inequalities: evidence from systematic reviews. *J Epidemiol Community Health* 2010;64:284–91.
- Papadopoulou SK, Hassapidou MN, Katsiki N, *et al.* Relationships between alcohol consumption, smoking status and food habits in Greek adolescents. vascular implications for the future. *Curr Vasc Pharmacol* 2017;15:167–73.
- Wong E. Tobacco use in king county. *Public Health Seattle & King County*, 2012. <https://www.kingcounty.gov/depts/health/data/~media/depts/health/data/documents/tobacco-use-in-king-county-may-2012.ashx>
- Bogan S, Donohue B. King County drug and alcohol deaths rose 9.5% in 2018. Available: <https://newsroom.uw.edu/news/king-county-drug-and-alcohol-deaths-rose-95-2018>
- Drug-caused deaths in King County, 2017. Available: <https://adai.washington.edu/WAdata/KingCountyDrugDeaths.htm>
- Gundlapalli AV, Carter ME, Palmer M, *et al.* Using natural language processing on the free text of clinical documents to screen for evidence of homelessness among US veterans. *AMIA Annu Symp Proc* 2013;2013:537–46.
- Savova GK, Masanz JJ, Ogren PV, *et al.* Mayo clinical text analysis and knowledge extraction system (cTAKES): architecture, component evaluation and applications. *J Am Med Inform Assoc* 2010;17:507–13.
- Gundlapalli AV, Carter ME, Divita G, *et al.* Extracting concepts related to homelessness from the free text of Va electronic medical records. *AMIA Annu Symp Proc* 2014;2014:589–98.
- Horng S, Sontag DA, Halpern Y, *et al.* Creating an automated trigger for sepsis clinical decision support at emergency department triage using machine learning. *PLoS One* 2017;12:e0174708.
- Feller DJ, Zucker J, Yin MT, *et al.* Using clinical notes and natural language processing for automated HIV risk assessment. *J Acquir Immune Defic Syndr* 2018;77:160–6.
- Dorr D, Bejan CA, Pizzimenti C, *et al.* Identifying patients with significant problems related to social determinants of health with natural language processing. *Stud Health Technol Inform* 2019;264:1456–7.



- 28 Berg K, Doktorchik C, Quan H. Meaningful information in the age of big data: a scoping review on social determinants of health data collection for electronic health records 2019.
- 29 2015 CDC HA-VTE prevention challenge champion. Available: <https://www.cdc.gov/ncbddd/dvt/documents/champ-fact-sheet-harborview.pdf>
- 30 Bulger EM, Kastl JG, Maier RV. The history of Harborview medical center and the Washington state trauma system. *Trauma Surg Acute Care Open* 2017;2:e000091.
- 31 Cronin RM, Fabbri D, Denny JC, *et al.* A comparison of rule-based and machine learning approaches for classifying patient portal messages. *Int J Med Inform* 2017;105:110–20.
- 32 Yao L, Mao C, Luo Y. Clinical text classification with rule-based features and knowledge-guided convolutional neural networks. *BMC Med Inform Decis Mak* 2019;19:71.
- 33 Medicare & Medicaid EHR Incentive Program. Meaningful use stage 1 requirements overview, 2010. Available: https://www.cms.gov/Regulations-and-Guidance/Legislation/EHRIncentivePrograms/downloads/MU_Stage1_ReqOverview.pdf
- 34 Quality Measures and Tobacco Cessation. Available: <https://www.bhthechange.org/wp-content/uploads/2017/12/Quality-Measures-and-Tobacco-Cessation.pdf>
- 35 Baclic O, Tunis M, Young K, *et al.* Challenges and opportunities for public health made possible by advances in natural language processing. *Can Commun Dis Rep* 2020;46:161–8.
- 36 Cai L, Zhu Y. The challenges of data quality and data quality assessment in the big data era. *Data Sci J* 2015;14:2.