

Harnessing information from injury narratives in the 'big data' era: understanding and applying machine learning for injury surveillance

Kirsten Vallmuur,¹ Helen R Marucci-Wellman,² Jennifer A Taylor,³ Mark Lehto,⁴ Helen L Corns,⁵ Gordon S Smith⁶

¹Queensland University of Technology, Centre for Accident Research and Road Safety—Queensland, Brisbane, Queensland, Australia

²Center for Injury Epidemiology, Liberty Mutual Research Institute for Safety, Hopkinton, Massachusetts, USA

³Department of Environmental & Occupational Health, School of Public Health, Drexel University, Philadelphia, Pennsylvania, USA

⁴School of Industrial Engineering, Purdue University, West Lafayette, Indiana, USA

⁵Center for Injury Epidemiology, Liberty Mutual Research Institute for Safety, Hopkinton, Massachusetts, USA

⁶National Center for Trauma and EMS, University of Maryland School of Medicine, Baltimore, Maryland, USA

Correspondence to

Dr Kirsten Vallmuur, Queensland University of Technology, Centre for Accident Research and Road Safety—Queensland, 130 Victoria Park Road, Kelvin Grove 4059, Brisbane, QLD 4053, Australia; k.vallmuur@qut.edu.au

Received 31 August 2015

Revised 1 December 2015

Accepted 8 December 2015

ABSTRACT

Objective Vast amounts of injury narratives are collected daily and are available electronically in real time and have great potential for use in injury surveillance and evaluation. Machine learning algorithms have been developed to assist in identifying cases and classifying mechanisms leading to injury in a much timelier manner than is possible when relying on manual coding of narratives. The aim of this paper is to describe the background, growth, value, challenges and future directions of machine learning as applied to injury surveillance.

Methods This paper reviews key aspects of machine learning using injury narratives, providing a case study to demonstrate an application to an established human-machine learning approach.

Results The range of applications and utility of narrative text has increased greatly with advancements in computing techniques over time. Practical and feasible methods exist for semiautomatic classification of injury narratives which are accurate, efficient and meaningful. The human-machine learning approach described in the case study achieved high sensitivity and PPV and reduced the need for human coding to less than a third of cases in one large occupational injury database.

Conclusions The last 20 years have seen a dramatic change in the potential for technological advancements in injury surveillance. Machine learning of 'big injury narrative data' opens up many possibilities for expanded sources of data which can provide more comprehensive, ongoing and timely surveillance to inform future injury prevention policy and practice.

INTRODUCTION

Injury narratives have long been recognised as valuable sources of information to understand injury circumstances and are increasingly available in the era of 'big data'. Narrative text mining and machine learning techniques have been developed that take advantage of greatly increased computing power and 'big data' to make predictions based on algorithms constructed from the data. However, along with the opportunities, challenges in adequately accessing and using injury narratives for public health surveillance and prevention exist. In this paper the authors describe the background, growth and utility of machine learning of injury narratives. A case study is also provided to demonstrate the application of an established human-machine learning approach. The authors then discuss the challenges and future directions of machine learning as applied to injury surveillance.

BACKGROUND

The 1990s marked the beginning of the electronic era; email and the internet were surfacing and electronic records took the form of .dbf files transcribed from hard copy files. In a 1997 article Sorock *et al* identified innovative approaches to improvements in work-related injury surveillance that reflected the utility of electronic records at this time.¹ These include¹ the use of narrative text fields from injury databases to extract useful epidemiological data,² data set linkage for aiding in incidence rate calculations and³ the development of comprehensive company-wide injury surveillance systems. Now almost 20 years later, the opportunities have expanded greatly; large amounts of coded injury data and text descriptions of injury circumstances (injury narratives) are being collected daily and are available in real time. However, while there have been some collective efforts to standardise injury data collection and classification systems, very little has been done to develop and standardise machine learning approaches using injury narratives.

WHO guidelines specify the following requirements for injury surveillance: to facilitate ongoing data collection, in a systematic way, which enables analysis and interpretation for timely dissemination which can be applied to prevention and control.² However, often injury information (for morbidity and mortality incidence reporting) is collected and may be classified without considering these requirements. While the data may be coded according to a standardised classification protocol (eg, International Classification of Diseases (ICD) coding in hospitals) the people assigning the codes are often administrative staff classifying the cases for billing purposes (not for prevention), with little professional training, although hospital discharge data is usually coded by a professional nosologist. In order to get these data recoded in such a way as to satisfy the requirements of surveillance requires significant investment and resources.

On the other hand there are some national agencies such as the National Center for Health Statistics which in addition to mortality coding use their nosologists to classify medical conditions, drugs and injuries reported in their large national health surveys in the USA (eg, the National Health and Nutrition Examination Survey and the National Health Survey). Coding systems useful to injury epidemiologists include: the ICD, International Classification of External Causes,³

To cite: Vallmuur K, Marucci-Wellman HR, Taylor JA, *et al.* *Inj Prev* Published Online First: [please include Day Month Year] doi:10.1136/injuryprev-2015-041813

and Nordic Classification of External Causes.⁴ Occupational injury surveillance systems however usually assign and use separate coding strategies aimed at identifying work exposures such as the National Institute for Occupational Safety and Health (NIOSH) Occupational Injury and Illness Classification System (OIICS)⁵ and the Type Of Occurrence Classification Scheme.⁶ These codes are often used for surveillance. However, even if the time and resources have been allotted to having trained coders assign these codes, there are still limitations in using the coded data alone. These include the limited scope, breadth and depth of injury mechanisms and scenarios captured from the codes (specifically reducing their value for injury prevention and control) and reliance on predetermined circumstances that may not capture all or the very unique case scenarios,⁷ nor all relevant injury factors (host, agent, vector, environment) contributing to an injury event as defined by Haddon.⁸

The utility of injury narratives for surveillance

Two recent reviews^{9 10} outlined a range of benefits for using narratives as a supplement to the restrictions of coded data, including: the identification of cases not able to be detected from coded data elements alone, extracting more specific information than codes allow, extracting data fields which aren't part of the prior coding schemas, establishing chain of events, identifying causes without specific codes and assessing coding accuracy.

Narrative text analyses also enables the identification of rare or emerging events usually not found using administratively assigned codes, a critical concern in injury surveillance.^{11–14} Incident narratives in their raw form can also be available in a more 'timely' manner than coded data and are now being used in novel applications such as syndromic surveillance.^{15 16}

The range of applications and utility of narrative text has also increased with recent advancements in computing techniques. However, some of the earliest applications predate the ability to search text electronically and were simply to identify cases to overcome the lack of reported or coded data. These include using newspaper clipping services where people were paid to read newspapers and identify articles that reference any of the injury or fatality topics on a list related to clients' interests who had paid the service to look for articles containing target words about specific companies.^{17 18} Now that news articles are on the web, computerised search has greatly simplified the process of searching for injury incidents using services such as Nexus.

Nowadays, with significant increases in the technological capabilities and capacity of computer systems, injury narratives which contain essential information about how the injury event occurred are more widely available in an 'ongoing' manner across a range of agencies (including but not limited to emergency services/first responders (ambulance, fire service, police), emergency departments/hospitals/trauma registries, coronial systems, occupational health and safety, insurance/compensation agencies (workplace/health/motor vehicle), consumer safety agencies, news services and even social networking sites (twitter/facebook), etc).

However, using these data for surveillance has historically proven cost-prohibitive and fraught with human error. Bertke *et al*¹⁹ reported that it took a single researcher 10 h (over the course of a few weeks to mitigate fatigue) to code 2400 workers' compensation injuries. Taylor *et al*²⁰ reported 100 total hours for three coders to discern cause of injury and reconcile differences from firefighter near-miss and injury narratives. As a database grows, the additional resources required to code the records become increasingly labour-prohibitive,

cost-prohibitive and time prohibitive. Only recently has the use of computerised coding algorithms enabled large-scale analysis of narrative text, presenting an efficient and plausible way for individuals to code large narrative data sets with accuracies of up to 90%.^{19 21} While autocoding increases accuracy and efficiency, it does not eliminate the need for human review entirely as humans must initially train the algorithm and conduct post hoc quality review.

There have been some limited situations where automated classification of injury narratives has become integrated into routine processes for national statistical purposes to reduce the amount and costs of manual coding, improve coding uniformity and reduce the time taken to process records. For example, many countries use software to automatically process injury text recorded on death certificates for broad ICD cause of death coding²² and the NIOSH in the USA has made available an online tool to aid state public health organisations in determining NIOSH occupation and industry codes.²³ These software programs built over several decades allow a substantial subset of records to be automatically coded usually with the caveat of limited accuracy. The accuracy however can often be improved if the algorithm is able to identify those which would be more accurately coded by humans (or should be unclassifiable) or that the software cannot confidently assign a code.

Over the past two decades, several authors of this paper have completed a number of studies^{1 24 25 21 26 27 20} on the utilisation of computer algorithms to streamline the classification of the event (or causes) documented in injury narratives for surveillance purposes. Their focus has been to create machine learning techniques to quickly filter through hundreds of thousands of narratives to accurately and consistently classify and track high-magnitude, high-risk and emerging causes of injury, information which can be used to guide the development of interventions for prevention of future injury incidents.²⁸ The results of this work have enabled the annual classification of very large batches of workers compensation (WC) claim incident narratives into Bureau of Labor Statistics (BLS) OIICS event codes for input in deriving the annual Liberty Mutual Workplace Safety Index—a surveillance metric ranking the leading causes (in terms of direct WC cost) of the most disabling work-related injuries in the USA every year.²⁹

Table 1 also provides examples of other studies, describing early uses and other more complex uses of narrative text. These examples include the integration of machine learning techniques to demonstrate the changing nature of this field.

CASE STUDY

To demonstrate one successful approach to the use of machine learning to classify injury narratives, the following case study briefly summarises a recent study by Marucci-Wellman *et al*²⁶ that accurately classified 30 000 WC narratives into injury events using a human-machine learning approach in order to match cost of claims by event category with national counts from the BLS Survey of Occupational Injury and Illness data. Coders who had been trained extensively on the BLS OIICS read each claim incident narrative on the case and classified the event that led to work-related injury into 1 of approximately 40 two-digit event codes. The data set was divided into a training set of 15 000 cases used for model development, and a prediction data set of 15 000 cases used for evaluating the algorithms performance on new narratives. A sample of WC claims incident narratives with BLS OIICS code assignments are shown below:

Table 1 Examples of original and complex applications of narrative text over time

Original applications			More complex applications		
Article details	Technique	Application	Article details	Technique	Application
Archer <i>et al</i> ¹⁸	Newspaper clipping service used to manually identify cases of firearm-related injuries (unintentional and intentional) along with other sources obtained from hospital, police and vital statistics	Newspaper clipping service identified almost a third of firearm-related cases (but only 17% of suicides) and were a cheap, accessible and simple data source albeit incomplete, especially for suicide	Homan <i>et al</i> ³⁰	Extracted 200 tweets from 2.5 million tweets which noted suicide terms, used expert and novice ranks of tweets for distress levels, and used support vector machine approach to topic model data	Automated tweet classification by distress levels to enable identification of individuals at risk of suicide through social media, with use of expert coders for training data and machine learning model choice, important factors affecting performance of model
Hume <i>et al</i> ³¹	Free text search of emergency department data from one New Zealand hospital for 1 year for one product (trampoline)	Identified the number of trampoline-related incidents and allowed case identification to enable further review of text and manual coding of extra circumstance details	Chen <i>et al</i> ³²	Automatic classification of mechanism and object categories for 15 000 emergency department cases across multiple hospitals using machine learning (matrix factorisation approach)	Classified mechanism and objects quickly with accuracy of 0.93, showing potential for use to reduce need for manual coding for injury surveillance, though need for expert input into modelling required throughout process to improve algorithm performance
Sorock <i>et al</i> ³³ and Lehto and Sorock ²⁴	Free text search of motor vehicle insurance claims database for 4 years to identify claims where road work occurring and key word categorisation of precrash activities and crash types through word frequency count and manual grouping of similar words to prepare keyword search strategy. Expanded to test a Bayesian modelling approach in second paper	First paper identified number of incidents and categorised precrash activities and crash types to examine patterns of incidents. Second paper established Bayesian approach more accurately classified cases than keywords and pointed to the early potential for Bayesian approaches to be developed in this field	Taylor <i>et al</i> ²⁰	Classified 2285 firefighter occupation-specific narratives (longer narratives with average of 216 words), with near-misses and injuries into injury mechanism and injury outcome using fuzzy and naïve Bayesian models with single word predictors	Classified external causes with accuracy of 0.74 using fuzzy model and 0.678 using Naïve model, with increased training set size producing higher sensitivity. Showed that Bayesian methods can be used for coding long narratives for injury incidents and near-misses
Bauer and Sector ³⁴	Development of a keyword based search to identify extent of product involvement in injury from emergency department based injury surveillance database, as well as use of an expert panel to assess preventability and potential for product safety responses	Ability to flag cases where high likelihood of consumer product involvement (defective, maladapted or intrinsically risky) and identify products most commonly associated with each category	Pan <i>et al</i> ³⁵	Use of named entity recognition techniques to automatically parse unstructured data from a range of databases (including RAPEX (Rapid Exchange of Information system) CPSC and product safety databases in China and Japan). Used Bayesian network approach to identify and code safety factors pertaining to electric shock	Automated extraction and coding of relevant cases incorporating a number of large publically available databases from different regions. Identification of the key safety factors involved in electric shock incidents (near miss and injuries), showing potential of multiple databases to extract common scenarios
Bondy <i>et al</i> ³⁶	Manual review of 4000 injury text reports from construction of Denver International Airport, and expert classification of case details according to Haddon's Matrix framework	Classification of text reports according to Haddon's Matrix framework provided a more complete injury description than only coding certain injury elements, as well as providing richer data to understand injury scenario and target prevention activity	Zhao <i>et al</i> ³⁷ and Zhao <i>et al</i> ³⁸	Use of electrocution text reports in national occupational injury database to extract either key features according to hierarchy of control framework or Haddon's Matrix framework. Used narrative text analysis (such as word clusters, entity extraction, word tagging and 'textual tag clouds') using NVivo qualitative software	Automated extraction and tagging of key features of reports and grouping according to overarching injury prevention frameworks, to examine main prevention foci as well as illustrate decision making chains. Demonstrates the utility of text analysis to extract and elucidate more complex injury causation scenarios

1. "STANDING UP FROM BENDING OVER STRUCK BACK ON MAID CART" ->Classified as BLS OIICS event code 63—struck against object or equipment.
2. "FELT PAIN WHILE PULLING LOAD OF WOOD WITH PALLET JACK" ->Classified as BLS OIICS event code 71—overexertion involving outside sources.
3. "STOPPED AT STOP SIGN WHEN REAR-ENDED BY ANOTHER VEH." ->Classified as BLS OIICS event code 26—Roadway incidents involving motorised land vehicle.
4. "SLIPPED AND FELL ON UNK SURFACE TWISTING HIS ANKLE SPRAININGIT".->Classified as BLS OIICS event code 42—Falls on same level.
5. "EMPLOYEE WAS WALKING ON THE STREET WHEN HIS RIGHT KNEE POPPED" ->Classified as BLS OIICS event code 73—Other exertions or bodily reactions.
Using the 15 000 narratives and manually assigned codes from the training set, a keyword list was created by parsing the words in each narrative (eg, standing, up, from, bending, etc). The occurrence or probability of each word in each category (P_{n_j/C_i}) was calculated as well as the marginal probability of each event category in the training data set ($P(C_i)$); these are the two parameters necessary for the reduced Naïve Bayes algorithm.²⁶ These statistics calculated from the training narratives were stored in a probability table and used to train the

algorithm. A similar word list and probability table was constructed for 2, 3 and 4 word sequences (each sequence considered as a keyword, eg, standing-up, up-from, from-bending, standing-up-from, etc). The Naïve Bayes model was used to assign a probability to each event code based on the keywords present in a particular narrative. The event code with the largest estimated probability was then chosen as the prediction for the words present.

The theoretical basis for the Naïve Bayes classifier and detailed instructions on how to implement the algorithm with narrative data have been thoroughly defined previously.^{21–26} Various software packages are now publically available for training (or building) the models based on the training data set and then making subsequent predictions. Weka³⁹ and Python⁴⁰ are two examples of publically available, easily downloadable and easily adaptable packages for development of the Naïve Bayes model. For this study, the Textminer software developed by one of the authors (ML) was used. The narratives were used in their raw form; although improved performance can be expected when misspellings are cleaned and words that have the same meaning are morphed into one syntax, the aim was to demonstrate what could be achieved by machine learning with little preprocessing of the narratives. However, a small list of frequently occurring ‘stop words’ believed to have little meaning for the classification assignment (eg, a, and, left, right) was removed from the narratives prior to calculating probabilities.

Two Naïve Bayes algorithms were run on each of the 15 000 prediction narratives using first the set of single keyword probabilities and second the sequenced keyword probabilities (stored in probability tables) from the training narratives in order to assign two independent computer generated classifications to the 15 000 prediction narratives.

The authors²⁶ found while the overall sensitivity of the two independent models was fairly good (0.67 naive_{sw}, 0.65 naive_{seq}), both algorithms independently predicted some categories much better than others, skewing the final distribution of the coded data (χ^2 $p < 0.0001$), and most of the cases in the smaller categories were not found. The sequence-word model showed improved performance where word order was important for differentiating causality. Still many categories had low performance. We consequently integrated a rule where we would *only* use the computer classifications when the two models agreed and then would manually code the remaining narratives. Implementing this rule resulted in an overall sensitivity of codes for the final coded data set of 87% with high sensitivity and PPVs across all categories (see [tables 2 and 3](#) and Marucci-Wellman *et al*²⁶ for more details). Note, high sensitivity and PPV are important for resulting in a final unbiased distribution of the coded data for surveillance and targeting prevention efforts. Also using this human-machine pairing resulted in 68% of the narratives coded by the algorithm, leaving only 32% to be coded by humans.

The authors found the accuracy of the human-machine system was at least as good and likely was even better than manual coding alone of all 15 000 records as the system uses consistent rules. This was demonstrated by comparing the results with inter-rater reliability data for four well trained human coders. While the evaluation of inter-rater reliability relies on different metrics, the inter-rater reliability performance of the four coders does not appear to be as systematically high and consistent as what is projected from the sensitivity and PPV values of the human-machine pairing method for the very large categories, nor the very small categories. Other readily available and easily adaptable machine learning techniques for narrative text analyses other than the Bayesian algorithms exist such as

support vector machine (SVM) and logistic regression (LR) and could also be incorporated to improve accuracy. Work has begun to investigate ensembles consisting of agreement between these various algorithms with some slightly improved results over the ones presented in the case study summary (see [table 4](#)). Overall, this case study demonstrates that a practical and feasible method exists for human-machine learning of short injury narratives. The computer was able to accurately classify many of the narratives of a large WC data set leaving a third for human review and resulting in a very high overall accuracy and very high accuracy across almost all categories (large and small) in the final coded data set. Accuracy can be further improved when a per cent of difficult cases, predicted by the algorithm with a low confidence, are rejected for manual coding.

DISCUSSION: CHALLENGES AND FUTURE DIRECTIONS

As illustrated in the previous case study, the use of off-the-shelf machine learning methods combined with human review of weakly predicted cases is an effective, easily applied method. However, this approach still required developing a large training set of previously coded cases to develop the model and then subsequent human review of around a third of the cases to attain high sensitivities across all categories in the prediction set. In practice, obtaining a good training set and the need for human review (which could be substantial if a third of a very large data set still requires manual coding) may be major application bottlenecks. Numerous strategies and approaches for tailoring methods to address this problem exist. For the most part, these strategies and approaches can be roughly divided as: focusing on obtaining more data (a larger training set), applying better learning algorithms, or going beyond the training set, using other sources of information, causal models, or human knowledge to preprocess the information used by the learning algorithm. The following discussion briefly builds on ideas generated by the case study and introduces some of these other approaches, their effectiveness and emerging trends in their use.

Obtaining more data or applying better algorithms

The use of a larger training set and better learning algorithms are commonly suggested strategies for improving model predictions. Previous work³² has shown that model performance improves for short injury narratives with larger training sets. The latter study also showed that SVM algorithms performed better than Naïve Bayes and several other learning algorithms. However, the improvements were clearly slowing down as the increase of training data continued. Furthermore, smaller categories were often poorly predicted by the algorithm, just as found in the case study above for Naïve Bayes, LR and SVM. Some further improvements in the SVM model performance were also observed by Chen *et al*³² after model factorisation using singular value decomposition (SVD) to map the word vectors to a lower dimensional space. The latter result was consistent with earlier studies showing improvements after feature space reduction using SVD,^{41–42} and SVD approaches are likely to be especially useful in ‘big data’ applications where there is substantial training data available for mapping the lower dimensional space.

Preprocessing data

Overall though, the results using thousands of training examples across multiple studies suggest that it is doubtful that the need for human review will be completely eliminated with more data or by better learning algorithms alone for complex multiclass

Table 2 The accuracy of the human-machine classification system: implementation of a strategic filter† based on agreement between two Naïve Bayes algorithms

BLS OIICS two-digit event code	Gold standard‡ (n)	Human-machine system coding of all narratives§						% Agreement between two manual coders¶	Fleiss κ**
		n _{pred} ††	% _{pred} ‡‡,§§	Sen¶¶	95% CI	PPV***	95% CI		
1* Violence and other injuries by persons or animals									
11 Intentional injury by person	159	132	0.9	0.81	0.75 to 0.87	0.98	0.95 to 1.00	81–97	0.85
2* Transportation incidents									
24 Pedestrian vehicular incidents	120	117	0.8	0.78	0.71 to 0.86	0.80	0.73 to 0.88	57–78	0.65
26 Roadway incidents motorised land vehicle	650	672	4.5	0.98	0.97 to 0.99	0.95	0.93 to 0.97	93–96	0.94
27 Non-roadway incidents motorised land vehicle	136	122	0.8	0.80	0.73 to 0.87	0.89	0.84 to 0.95	52–84	0.62
4* Falls, slips, trips									
41 Slip or trip without fall	806	658	4.4	0.70	0.67 to 0.73	0.86	0.83 to 0.89	66–89	0.71
42 Falls on same level	2148	2386	15.9	0.92	0.91 to 0.93	0.83	0.81 to 0.84	85–93	0.86
43 Falls to lower level	1065	1176	7.8	0.89	0.87 to 0.91	0.81	0.79 to 0.83	78–92	0.81
5* Exposure to harmful substances or environments									
53 Exposure to temperature extremes	141	130	0.9	0.86	0.8 to 0.92	0.93	0.89 to 0.97	82–98	0.88
55 Exposure to other harmful substances	175	165	1.1	0.83	0.77 to 0.88	0.88	0.83 to 0.93	81–96	0.87
6* Contact with objects and equipment									
62 Struck by object or equipment	1651	1749	11.7	0.90	0.89 to 0.92	0.85	0.83 to 0.87	82–90	0.82
63 Struck against object or equipment	466	397	2.6	0.74	0.7 to 0.78	0.87	0.84 to 0.91	66–83	0.68
64 Caught in or compressed by equipment	505	532	3.5	0.90	0.87 to 0.93	0.86	0.83 to 0.89	72–83	0.75
7* Overexertion and bodily reaction									
70 Overexertion and bodily reaction, uns	188	151	1.0	0.59	0.51 to 0.66	0.73	0.66 to 0.80	6–48	0.19
71 Overexertion involving outside sources	4189	4334	28.9	0.95	0.95 to 0.96	0.92	0.91 to 0.93	87–95	0.87
72 Repetitive motions involving micro tasks	484	537	3.6	0.90	0.87 to 0.92	0.81	0.77 to 0.84	71–83	0.75
73 Other exertions or bodily reactions	916	827	5.5	0.79	0.76 to 0.82	0.88	0.85 to 0.90	56–85	0.64
X* All other classifiables (n<100) in training data set									
xx Other small (n<100 cases) classifiable categories†††	632	467	3.1	0.68	0.64 to 0.72	0.92	0.89 to 0.94	–	–
Non-classifiable									
9999 Non-classifiable	569	448	3.0	0.70	0.66 to 0.74	0.89	0.86 to 0.92	69–84	0.72
Overall	15 000	15 000	100.0	0.87	0.87 to 0.88	0.87	0.87 to 0.88	77–90	0.78

Adapted from Marucci-Wellman *et al.*²⁶

†A filter is a technique to decide which narratives the computer should classify versus which should be left for a human to read and classify.

‡Gold standard codes were assigned to each narrative by expert manual coders.

§Human-machine system: The computer assigns codes to narratives that the algorithms agreed on the classification (68% of the data set), and the remainder are manually coded (32% of the data set).

¶¶Two-coder agreement, for example, 6 total comparisons, coder 1 compared with 2, 3, 4, coder 2 compared with 3, 4, coder 3 compared with 4.

**Fleiss κ between 0 and 1, >0.6 considered good agreement, >0.8 considered very good agreement.

††n_{pred}=number predicted into category.‡‡%_{pred}=per cent of cases in whole data set predicted into category.

§§The distribution of two-digit classifications will be skewed towards categories with high sensitivity, biasing the final distribution of the coded data sets.

¶¶¶Sen=Sensitivity: (true positives) the percentage of narratives that had been coded by the experts into each category that were also assigned correctly by the algorithm.

***PPV, the percentage of narratives correctly coded into a specific category out of all narratives placed into that category by the algorithm.

†††Two-digit categories with <100 cases.

BLS, Bureau of Labor Statistics; OIICS, Occupational Injury and Illness Classification System.

*-Asterisks denote a summary level code not assigned to individual cases (see http://www.bls.gov/iif/osh_oiccs_2010_2_4_2.pdf)

coding schemes and especially so when there is a need to assign rarely occurring categories (ie, needle stick injuries in the case study). One potentially promising strategy for improving performance for smaller categories is to go beyond the training set, using other sources of information, causal models or human knowledge, to preprocess the information used by the learning algorithm. Numerous approaches have been used for preprocessing injury text prior to applying the learning algorithms such as word stemming, lemmatisation, dropping infrequent or frequent

words, or weighting schemes such as TF-IDF.³² One advantage of such approaches is that they provide an easy way of reducing the dimensionality of the word vector, which can speed learning of any machine learning algorithm. However, this may sacrifice accuracy, with the authors' preliminary work using Naïve Bayes, LR and SVM showing that these preprocessing approaches have the potential to reduce the overall detection (distinguishing between categories) capability, and especially for small categories.⁴³ Part of the problem is that such approaches do not

Table 3 The accuracy of the human-machine classification system: implementation of a strategic filter† based on agreement between the two naïve bayes algorithms (results for small categories only, n<100 cases in each category)

BLS OIICS two-digit event code	Gold standard‡ (n)	Human-machine system coding of all narratives§					% Agreement between two manual coders¶	Fleiss κ** manual coders
		n _{pred} ††	Sen‡‡	(95% CI)	PPV§§	95% CI		
1* Violence and other injuries by persons or animals								
12 Injury by person—intentional or intent unknown	96	78	0.66	0.56 to 0.75	0.81	0.71 to 0.88	47–78	0.57
13 Animal and insect related incidents	99	79	0.80	0.71 to 0.87	1.00	1.00 to 1.00	79–94	0.87
2* Transportation incidents								
20 Transportation incident, unspecified	3	3	1.00	1.00 to 1.00	1.00	1.00 to 1.00	0–0	0.00
21 Aircraft incidents	22	15	0.68	0.47 to 0.89	1.00	1.00 to 1.00	0–75	0.17
22 Rail vehicle incidents	6	4	0.67	0.12 to 1.00	1.00	1.00 to 1.00	0–100	0.67
23 Animal and other non-motorised vehicle transport incidents	14	13	0.86	0.65 to 1.00	0.92	0.76 to 1.00	0–0	0.00
25 Water vehicle incidents	11	5	0.45	0.1 to 0.81	1.00	1.00 to 1.00	0–88	0.25
3* Fires and explosion								
31 Fires	22	20	0.91	0.78 to 1.00	1.00	1.00 to 1.00	55–88	0.58
32 Explosions	21	18	0.86	0.69 to 1.00	1.00	1.00 to 1.00	44–83	0.46
4* Falls, slips, trips								
40 Fall, slip, trip, unspecified	4	2	0.50	0.00 to 1.00	1.00	1.00 to 1.00	0–0	0.00
44 Jumps to lower level	57	39	0.61	0.48 to 0.74	0.90	0.80 to 1.00	51–90	0.65
45 Fall or jump curtailed by personal fall arrest system	3	2	0.67	0.00 to 1.00	1.00	1.00 to 1.00	0–0	0.00
5* Exposure to harmful substances or environments								
50 Exposure to harmful substances or environ, unspecified	23	18	0.78	0.6 to 0.96	1.00	1.00 to 1.00	21–88	0.33
51 Exposure to electricity	27	18	0.67	0.48 to 0.86	1.00	1.00 to 1.00	65–88	0.81
52 Exposure to radiation and noise	38	36	0.87	0.76 to 0.98	0.92	0.82 to 1.00	54–100	0.80
54 Exposure to air and water pressure change	1	0	0.00	–	0.00	–	0–100	0.40
57 Exposure to traumatic or stressful even nec	32	23	0.72	0.55 to 0.88	1.00	1.00 to 1.00	73–85	0.80
59 Exposure to harmful substances or environments, nec	1	7	0.00	–	0.00	–	0–100	0.12
6* Contact with objects and equipment								
60 Contact with objects and equipment, uns	78	43	0.54	0.43 to 0.65	0.98	0.93 to 1.00	12–63	0.25
61 Needle stick	1	1	1.00	1.00 to 1.00	1.00	1.00 to 1.00	–	–
65 Struck/caught/crush in collapsing structure, equip or material	5	3	0.60	0.00 to 1.00	1.00	1.00 to 1.00	0–0	0.33
66 Rubbed or abraded by friction or pressure	16	12	0.69	0.43 to 0.94	0.92	0.73 to 1.00	0–50	0.11
67 Rubbed abraded or jarred by vibration	7	4	0.57	0.08 to 1.00	1.00	1.00 to 1.00	0–67	0.14
69 Contact with objects and equipment, nec	1	1	1.00	1.00 to 1.00	1.00	1.00 to 1.00	–	–
7* Overexertion and bodily reaction								
74 Bodily conditions nec	20	10	0.50	0.26 to 0.74	1.00	1.00 to 1.00	0–75	0.33
78 Multiple types of overexertion and bodily reactions	23	13	0.39	0.18 to 0.61	0.69	0.40 to 0.98	0–0	0.00
79 Overexertion and bodily reaction and exertion, nec	1		0.00	–	0.00	–	–	–
Overall	437	467	0.68	0.64 to 0.72	0.92	0.89 to 0.94		

Adapted from Marucci-Wellman *et al.*²⁶

†A filter is a technique to decide which narratives the computer should classify versus which should be left for a human to read and classify.

‡Gold standard codes were assigned to each narrative by expert manual coders.

§Human-machine system consisted of human coding 32% of the data set, machine coding 68% of the data set.

¶Two-coder agreement, for example, 6 total comparisons, coder 1 compared with 2, 3, 4, coder 2 compared with 3, 4, coder 3 compared with 4.

**Fleiss κ between 0 and 1, >0.6 considered good agreement, >0.8 considered very good agreement.

††n_{pred}=number predicted into category.

‡‡Sen=Sensitivity: (true positives) the percentage of narratives that had been coded by the experts into each category that were also assigned correctly by the algorithm.

§§PPV, the percentage of narratives correctly coded into a specific category out of all narratives placed into that category by the algorithm.

BLS, Bureau of Labor Statistics; OIICS, Occupational Injury and Illness Classification System.

*-Asterisks denote a summary level code not assigned to individual cases (see http://www.bls.gov/iif/osh_oicis_2010_2_4_2.pdf)

Table 4 The accuracy of the human-machine classification system: implementation of a strategic filter* based on agreement of predictions between selected combinations of different algorithms (Naïve Bayes single word, Naïve Bayes bi-gram, SVM, logistic regression)

Models	Two-model agreement				Three-model agreement		
	SVM=Naïve Bayes single word (%)	SVM=Naïve Bayes bi-gram (%)	SVM=Logistic (%)	Logistic=Naïve Bayes single word (%)	Logistic=Naïve Bayes bi-gram (%)	SVM=Naïve Bayes single word=logistic (%)	SVM=Naïve Bayes single word=Naïve Bayes bi-gram (%)
Overall Sensitivity/PPV	87	89	81	86	88	89	93
Manually coded	28	33	14	24	29	31	43

*A filter is a technique to decide which narratives the computer should classify versus which should be left for a human to read and classify. SVM, support vector machine.

consider the meaning of words. For example, in related as yet unpublished work, the authors found that stemming or lemmatising the words ‘lifting’ and ‘lifts’ to their root ‘lift’ reduces the ability of SVM, NB and LR to distinguish injuries related to exertion from those caused by man lifts or fork lifts. Similarly, dropping infrequent words in this large word set of 10 000 words such as ‘muggers’ or ‘rape’ reduced the ability to identify assault cases.

Targeted mapping of only certain words to a common meaning, on the other hand, tended to improve performance (eg, HOT and SCALDING or bike and bicycle). The latter approach was especially useful for finding predictive word sequences (eg, ‘all words that mean a person’ followed by the word ‘fell’ separates struck by events from fall events). Based on the author’s preliminary results, systematic development of a lexicon mapping words, word-sequences and word combinations that relate to important concepts can greatly improve the sensitivity across categories of any machine learning algorithm. For example, the authors found the generic concept ‘hit body part on’ identified as a sequence of words that can mean hit, followed by words that can mean a body part, followed by either the frequent words ‘or’ or ‘against’, greatly improved the ability of Naïve Bayes, SVM and LR alike to distinguish struck against events from falls and struck against events. The finding that a good lexicon can improve the performance of machine learning algorithms for short injury narratives is not surprising. The caveat is that manually developing a good lexicon is very time-consuming, since data sets will contain thousands of unique words and words will have different meanings depending on what other words are present (really requiring topic appropriate linguist experts to do this work). Further complicating the matter, a causal model may be necessary to organise the concepts into a predictive model. Illustrating recent developments in this direction, Abdat *et al*⁴⁴ developed a causal model of construction injuries using a Bayesian network to identify the probable explanation of injuries based on generic factors extracted by experts from injury scenarios. Other work in this direction included the use of automated named entity recognition techniques to automatically parse unstructured data from several databases which were then used in a Bayesian network to identify and code safety factors.³⁵

An interesting conjecture is that these findings suggest a lexicon or causal factors generated from one text mining project can be used to help code another project’s uncoded narratives. Transfer of results would seem to be especially promising when data sets have the same focus, like occupational hazards. For example, if the results obtained using the database from the National Firefighter Near-Miss Reporting System²⁰ were applied to narratives from the Fire Fighter Fatality Investigation and Prevention Program, one would expect falls to be predicted with fairly good accuracy because

the language firefighters use to describe their hazards is similar (‘roof, spongy’ are precise predictors for firefighter falls caused from a weakening roof on fire). Similarly, a multitude of terms identified as toxic chemicals (eg, hydrogen sulfide, toluene) in one data set could be directly mapped to the concept ‘toxic chemical’ used in a new application, rather than relying on the training set alone. Future studies might also explore how well key words and word predictors in a home and leisure injury database²⁵ would predict injuries in occupational narratives. If one wanted to autocode causes of injury in firefighter narratives using results obtained from a knowledge database (meaning a collection of either narratives linked to manually assigned codes or word lists with corresponding probability weights) created from a home and leisure population level database, the terms used to describe important concepts in a firefighter database could be nodes in a Bayesian network retrained using the home and leisure injury database to estimate probability weights (P_n/C_i) for the new database. The new weights would adjust the original weights for terms such as ‘roof, spongy’ used as a precise predictor for firefighter falls but unlikely to indicate a fall when at home or in leisure activities. This approach will enable the development of weighting coefficients (as adjustments) to the probabilities that comprise the knowledge database before it is transferred from population narratives to occupational narratives. This work—while currently hypothetical—would, if feasible, provide critical proof of concept: if high specificity, sensitivity and PPV are able to be attained, there would be good evidence that weighting of probabilities would be the next step in making machine learning algorithms more broadly transferrable helping to reduce resources needed for human coding.

Building an open source knowledge base

For machine learning algorithms to be broadly used, they need to be accessible and refined in an open source manner. Ideally, researchers could share data and algorithms, perhaps in a cloud based shared-access knowledge database. Along these lines, Purdue University (ML) is in the process of creating an open source framework that can serve as a repository for shared injury coding knowledge databases. This framework would allow remote access to data sets of coded and uncoded narratives, machine learning algorithms, lexicons and other information, enabling researchers to share their results, develop better models more quickly and ultimately reduce the need to manually code in the traditionally resource-dependent manner. The expectation is that as the open source repository grows, new models will be developed that accurately code injury narratives within specific content areas. As more narratives are put into the knowledge database such models should perform more precisely and accurately. The end product would be an open-sourced knowledge

repository that stores words and associated probabilities in order to code injury narratives, where researchers and other organisations may upload their injury narratives, select what rubric and algorithm to apply, and then run the model to obtain injury codes for their narrative data.

Providing better access to training data and cloud based computer coding methods would enable researchers without previous access to computerised coding software and/or without a training set for the algorithm to code their data. This has global implications because health systems in the developing world have yet to move to computerised information systems and their only option may be narratives as trained coders are often scarce.

A shared knowledge database would enable injury researchers, organisations and government health agencies to code and analyse large injury narrative data sets without the need for substantial resources as previously required, liberating these untapped data sources to be used for surveillance, policy and implementing interventions. Ultimately, the future of injury surveillance must address who funds such a data warehouse and how it is financially sustained with appropriate technical assistance.

One of the challenges in building a knowledgebase of narratives and moving from privately used data sets to publically available data sets is the issue of confidentiality. Injury narratives may contain personally identifiable information (such as patient names) or company identifiable information (such as brands of products). To enable sharing of narratives more publically, language parsing techniques which can automatically de-identify details from narrative text (without losing the context of the narrative) will need to be incorporated into text mining methods, and there have already been significant advances in such techniques.⁴⁵

Human-directed learning

Nevertheless, algorithms do only what humans tell them. The human factors of manual review, quality assurance, and 'knowing your data' will still be required especially to identify new or emerging hazards and to understand the complex interaction of contributory factors—a principle of surveillance. Text mining for injury surveillance stands apart from other data mining efforts such as that used by generic search engines. Generic search engines allow algorithms to find whatever they can, while human-directed injury surveillance through text mining is looking for *particular* outcomes—injuries, and particular features (eg, host, agent, vector environment), classifiable to specified categories defined by the end user. The role of the human in teaching the algorithm how to behave is vital to getting it right.

It is difficult for an algorithm on its own to be able to assign classifications in all categories with the same level of confidence and very difficult to improve the accuracy of computer generated codes for the small categories or for identifying emerging hazards. Improvement beyond simply modelling of a training data set to use on a prediction data set requires either sophisticated filtering or tailoring of the algorithm (with natural language processing) to identify small categories or other nuances of the coding protocol and the latter approach will still not allow for emerging risks to surface.

It was stated from the beginning²⁵ that manual coding should never be completely replaced and therefore a best practice approach should incorporate some manual coding, assigning a computer classification only for more repetitive events where the models are able to confidently predict the correct classification. This will be especially important for rare events and/or

emerging hazards that appear only a very small number of times or not at all in a training data set. For example, a new MVC hazard (exploding magnesium steering column) would cause a human reviewer to query why steering columns explode on impact and if they represent a new material hazard to drivers and first responders. An algorithm would simply say this does not happen enough to be coded with certainty and would flag it for manual review. For large administrative data sets, incorporation of methods based on human-machine pairings such as presented in this paper using readily available off-the-shelf machine learning techniques result in only a fraction of narratives that require manual review.

CONCLUSION

Machine learning of 'big injury narrative data' opens up many possibilities for expanded sources of data that can provide more comprehensive, ongoing and timely surveillance to inform injury prevention policy and practice in the future. This paper has demonstrated the significant value that injury narratives provide beyond structured coded data sets. It is critically important that, as an injury prevention

What is already known on this subject

- ▶ Large amounts of coded injury data and injury narratives are being collected universally daily and are available real time, yet the development and standardisation of machine learning approaches using injury narratives is nascent.
- ▶ Injury narratives provide opportunities to (A) identify the cases not able to be detected due to coding limitations, (B) extract more specific information than codes allow, (C) extract data fields which aren't part of the coding schema, (D) establish chain-of-events scenarios and (E) assess coding accuracy.
- ▶ The main focus of machine learning techniques using injury narratives has been to quickly filter large numbers of narratives to accurately and consistently classify and track high-magnitude, high-risk and emerging causes of injury, to guide the development of interventions for prevention of future injury incidents.

What this study adds

- ▶ Reiteration of the significant value that injury narratives provide beyond structured coded data sets and evidence for the continued need to advocate for narratives to be included (or introduced) in routine data sources to capitalise on this potential as computing and technical capacity expands.
- ▶ Demonstration of a practical and feasible method for semiautomatic classification using human-machine learning of injury narratives which is accurate, efficient and meaningful and applicable to different injury domains.
- ▶ The opening of a dialogue within the injury surveillance community regarding future steps towards developing a 'big injury narrative data' knowledge base to allow for the building, testing and refinement of machine learning algorithms.

community, we continue to advocate for the need for narratives to be included (or introduced) in routine data sources to capitalise on this potential as computing and technical capacity expands and not just rely on coded checkboxes. Second, the authors have argued for the need for a more systematic and incremental approach to developing machine learning approaches for the specialised purpose of injury surveillance, as distinct from other applications of machine learning more broadly. Modelling techniques (and research applications) vary in terms of levels of specificity and sensitivity, simplicity and complexity, and the building and refinement of these techniques require input from content experts and technical experts. The authors proposed future steps towards developing a 'big injury narrative data' platform to allow for the building, testing and refinement of machine learning algorithms. Finally, the need for human-machine pairings was reiterated to ensure machine learning approaches continue to reflect the underlying principles of injury surveillance.

The last 20 years has seen a dramatic change in the potential for technological advancements in injury surveillance and we have many examples of successful applications of such technology to injury narratives. It is now time to consolidate these learnings to build more sustainable, reliable and efficient approaches which will ensure the most robust use of the evidence base for injury prevention.

Contributors KV planned the manuscript, drafted sections, consolidated and revised drafts from authors, and prepared the final manuscript. ML, HRM-W and HLC wrote the case study. JAT, ML, HRM-W, HLC and GSS provided drafts of sections and edited and revised consecutive drafts of the paper.

Funding KV is supported by an Australian Research Council Future Fellowship under Grant FT120100202. GSS is supported by a grant from the US National Institute on Alcohol Abuse and Alcoholism (R01AA18707).

Competing interests None declared.

Provenance and peer review Commissioned; externally peer reviewed.

REFERENCES

- Sorock GS, Smith GS, Reeve GR, *et al.* Three perspectives on work-related injury surveillance systems. *Am J Ind Med* 1997;32:116–28.
- World Health Organisation. *WHO Injury Surveillance Guidelines*. Geneva: World Health Organisation, 2001.
- World Health Organization (WHO). *International Classification of External Causes of Injury (ICECI)*. Geneva, 2003.
- Nordic Medico-Statistical Committee. *NOMESCO Classification of External Causes of Injuries (Fourth revised edition)*. Copenhagen: AN:sats, 2007.
- United States Department of Labor Bureau of Labor Statistics. *Occupational Injury and Illness Classification Manual, Version 2.01*. USA, 2012.
- Australian Safety and Compensation Council. *Type of Occurrence Classification System (TOOCS) Third Edition Revision Canberra*. Australia: Australian Government, 2008.
- McKenzie K, Fingerhut L, Walker S, *et al.* Classifying external causes of injury: history, current approaches, and future directions. *Epidemiol Rev* 2012;34:4–16.
- Runyan C. Introduction: back to the future—revisiting Haddon's conceptualization of injury epidemiology and prevention. *Epidemiol Rev* 2003;25:60–4.
- McKenzie K, Scott D, Campbell M, *et al.* The use of narrative text for injury surveillance research: a systematic review. *Accid Anal Prev* 2010;42:354–63.
- Vallmuur K. Machine learning approaches to analysing textual injury surveillance data: a systematic review. *Accid Anal Prev* 2015;79:41–9.
- Stout N. *Analysis of narrative text fields in occupational injury data*. Occupational injury. CRC Press; 1998.
- Bunn TL, Slavova S, Hall L. Narrative text analysis of Kentucky tractor fatality reports. *Accid Anal Prev* 2008;40:419–25.
- Lipscomb HJ, Glazner J, Bondy J, *et al.* Analysis of text from injury reports improves understanding of construction falls. *J Occup Environ Med* 2004;46:1166–73.
- Smith GS, Timmons RA, Lombardi DA, *et al.* Work-related ladder fall fractures: identification and diagnosis validation using narrative text. *Accid Anal Prev* 2006;38:973–80.
- Chapman WW, Christensen LM, Wagner MM, *et al.* Classifying free-text triage chief complaints into syndromic categories with natural language processing. *Artif Intell Med* 2005;33:31–40.
- Muscattello DJ, Churches T, Kaldor J, *et al.* An automated, broad-based, near real-time public health surveillance system using presentations to hospital Emergency Departments in New South Wales, Australia. *BMC Public Health* 2005;5:141.
- Rainey DY, Runyan CW. Newspapers: a source for injury surveillance? *Am J Public Health* 1992;82:745–6.
- Archer P, Mallonee S, Schmidt A, *et al.* Oklahoma firearm-related injury surveillance. *Am J Prev Med* 1998;15:83–91.
- Bertke S, Meyers A, Wurzelbacher S, *et al.* Development and evaluation of a naïve bayesian model for coding causation of workers' compensation claims. *J Safety Res* 2012;43:327–32.
- Taylor JA, Lacovara AV, Smith GS, *et al.* Near-miss narratives from the fire service: a bayesian analysis. *Accid Anal Prev* 2014;62:119–29.
- Lehto M, Marucci-Wellman H, Corns H. Bayesian methods: a useful tool for classifying injury narratives into cause groups. *Inj Prev* 2009;15:259–65.
- Ossiander E. Using Textual Cause-of-Death Data to Study Drug Poisoning Deaths. *Am J Epidemiol* 2014;179:884–94.
- Centers for Disease Control and Prevention. NIOSH Industry and Occupation Computerized Coding System (NIOCCS). 2015. <http://wwwn.cdc.gov/niosh-nioccs/>
- Lehto MR, Sorock GS. Machine learning of motor vehicle accident categories from narrative data. *Methods Inf Med* 1996;35:309–16.
- Marucci-Wellman H, Lehto MR, Sorock GS, *et al.* Computerized coding of injury narrative data from the National Health Interview Survey. *Accid Anal Prev* 2004;36:165–71.
- Marucci-Wellman HR, Lehto MR, Corns HL. A practical tool for public health surveillance: semi-automated coding of short injury narratives from large administrative databases using naïve bayes algorithms. *Accid Anal Prev* 2015;84:165–76.
- Marucci-Wellman H, Lehto M, Corns H. A combined fuzzy and naïve Bayesian strategy can be used to assign event codes to injury narratives. *Inj Prev* 2011;17:407–14.
- Horan JM, Mallonee S. Injury surveillance. *Epidemiol Rev* 2003;25:24–42.
- Marucci-Wellman HR., Courtney TK, Corns HL, *et al.* The direct cost burden of 13 years of disabling workplace injuries in the U.S. (1998–2010): Findings from the Liberty Mutual Workplace Safety Index. *J Safety Res* 2015;55: Published Online First: December 2015, 53–62.
- Homan C, Johar R, Liu T, *et al.* editors. Toward Macro-Insights for Suicide Prevention: Analyzing Fine-Grained Distress at Scale. *Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality Proceedings of the Workshop*; Baltimore, Maryland, USA; 2014.
- Hume PA, Chalmers DJ, Wilson BD. Trampoline injury in New Zealand: emergency care. *Br J Sports Med* 1996;30:327–30.
- Chen L, Vallmuur K, Nayak R. Injury narrative text classification using the factorization model. *BMC Med Inform Decis Mak* 2015;15(Suppl 1):S5.
- Sorock GS, Ranney TA, Lehto MR. Motor vehicle crashes in roadway construction workzones: an analysis using narrative text from insurance claims. *Accid Anal Prev* 1996;28:131–8.
- Bauer R, Sector M. Preventive product safety—monitoring accidental injuries related to consumer products in the European Union. *Inj Control Saf Promot* 2003;10:253–5.
- Pan S, Wang L, Wang K, *et al.* A knowledge engineering framework for identifying key impact factors from safety-related accident cases. *Syst Res and Behav Sci* 2014;31:383–97.
- Bondy J, Lipscomb H, Guarini K, *et al.* Methods for using narrative text from injury reports to identify factors contributing to construction injury. *Am J Ind Med* 2005;48:373–80.
- Zhao D, McCoy A, Kleiner B, *et al.* Control measures of electrical hazards: an analysis of construction industry. *Safety Sci* 2015;77:143–51.
- Zhao D, McCoy A, Kleiner B, *et al.* Decision-Making Chains in Electrical Safety for Construction Workers. *J Constr Eng Manage* 2016;142. doi:10.1061/(ASCE)CO.1943-7862.0001037, 04015055
- Hall M, Frank E, Holmes G, *et al.* The WEKA data mining software: an update. *SIGKDD Explorations* 2009;11:10.
- Pedregosa F. Scikit-learn: machine learning in python. *J of Mach Learn Res* 2011;12:2825–30.
- Noorinaeini A, Lehto M. Mathematical models of human text classification. In: Duffy V, eds. *Handbook of digital human modeling for human factors and ergonomics*. Mahwah, NJ: Lawrence Erlbaum Associates, Inc., 2009:17.1–17.5.
- Noorinaeini A, Lehto M. Hybrid singular value decomposition; a model of text classification. *Int J Hum Factors Model and Simulation* 2006;1:95–118.
- Huang H, Lehto M. Significance of low-frequency words in text classification of open-ended survey responses. *2nd Global Conference on Engineering and Technology Management*; 4–5 September 2015; Chicago, IL, USA, 2015.
- Abdat F, Leclercq S, Cuny X, *et al.* Extracting recurrent scenarios from narrative texts using a Bayesian network: application to serious occupational accidents with movement disturbance. *Accid Anal Prev* 2014;70:155–66.
- Deleger L, Molnar K, Savova G, *et al.* Large-scale evaluation of automated clinical note de-identification and its impact on information extraction. *J Am Med Inform Assoc* 2013;20:84–94.



Harnessing information from injury narratives in the 'big data' era: understanding and applying machine learning for injury surveillance

Kirsten Vallmuur, Helen R Marucci-Wellman, Jennifer A Taylor, Mark Lehto, Helen L Corns and Gordon S Smith

Inj Prev published online January 4, 2016

Updated information and services can be found at:

<http://injuryprevention.bmj.com/content/early/2016/01/04/injuryprev-2015-041813>

These include:

References

This article cites 32 articles, 8 of which you can access for free at:
<http://injuryprevention.bmj.com/content/early/2016/01/04/injuryprev-2015-041813#BIBL>

Email alerting service

Receive free email alerts when new articles cite this article. Sign up in the box at the top right corner of the online article.

Notes

To request permissions go to:

<http://group.bmj.com/group/rights-licensing/permissions>

To order reprints go to:

<http://journals.bmj.com/cgi/reprintform>

To subscribe to BMJ go to:

<http://group.bmj.com/subscribe/>