

Original Paper

Identifying Acute Low Back Pain Episodes in Primary Care Practice From Clinical Notes: Observational Study

Riccardo Miotto^{1,2,3}, PhD; Bethany L Percha^{2,3}, PhD; Benjamin S Glicksberg^{1,2,3}, PhD; Hao-Chih Lee^{2,3}, PhD; Lisanne Cruz⁴, MD, MSc, FAAPMR; Joel T Dudley^{2,3}, PhD; Ismail Nabeel⁵, MD, MPH

¹Hasso Plattner Institute for Digital Health at Mount Sinai, Icahn School of Medicine at Mount Sinai, New York, NY, United States

²Institute for Next Generation Healthcare, Icahn School of Medicine at Mount Sinai, New York, NY, United States

³Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, United States

⁴Department of Physical Medicine and Rehabilitation, Icahn School of Medicine at Mount Sinai, New York, NY, United States

⁵Department of Environmental Medicine and Public Health, Icahn School of Medicine at Mount Sinai, New York, NY, United States

Corresponding Author:

Ismail Nabeel, MD, MPH

Department of Environmental Medicine and Public Health

Icahn School of Medicine at Mount Sinai

17 East 102nd Street, Box 1043

New York, NY, 10029

United States

Phone: 1 (614) 423 9057

Email: ismail.nabeel@icahn.mssm.edu

Abstract

Background: Acute and chronic low back pain (LBP) are different conditions with different treatments. However, they are coded in electronic health records with the same International Classification of Diseases, 10th revision (ICD-10) code (M54.5) and can be differentiated only by retrospective chart reviews. This prevents an efficient definition of data-driven guidelines for billing and therapy recommendations, such as return-to-work options.

Objective: The objective of this study was to evaluate the feasibility of automatically distinguishing acute LBP episodes by analyzing free-text clinical notes.

Methods: We used a dataset of 17,409 clinical notes from different primary care practices; of these, 891 documents were manually annotated as *acute LBP* and 2973 were generally associated with LBP via the recorded ICD-10 code. We compared different supervised and unsupervised strategies for automated identification: keyword search, topic modeling, logistic regression with bag of n-grams and manual features, and deep learning (a convolutional neural network-based architecture [ConvNet]). We trained the supervised models using either manual annotations or ICD-10 codes as positive labels.

Results: ConvNet trained using manual annotations obtained the best results with an area under the receiver operating characteristic curve of 0.98 and an F score of 0.70. ConvNet's results were also robust to reduction of the number of manually annotated documents. In the absence of manual annotations, topic models performed better than methods trained using ICD-10 codes, which were unsatisfactory for identifying LBP acuity.

Conclusions: This study uses clinical notes to delineate a potential path toward systematic learning of therapeutic strategies, billing guidelines, and management options for acute LBP at the point of care.

(*JMIR Med Inform* 2020;8(2):e16878) doi: [10.2196/16878](https://doi.org/10.2196/16878)

KEYWORDS

electronic health records; clinical notes; low back pain; natural language processing; machine learning

Introduction

Low back pain (LBP) is one of the most common causes of disability in US adults younger than 45 years [1], with 10 to

20% of American workers reporting persistent back pain [2]. LBP impacts one's ability to work and affects the quality of life. For example, in 2015, Luckhaupt et al showed that, from a pool of 19,441 people, 16.9% of workers with any LBP and 19.0% of those with frequent and severe LBP missed at least

one full day of work over a period of 3 months [3]. LBP events also lead to a significant financial burden for both individuals and clinical facilities, with combined direct and indirect costs of treatment for musculoskeletal injuries and associated pain estimated to be approximately US \$213 billion annually [4].

LBP events fall into 2 major categories: acute and chronic [5]. Acute LBP occurs suddenly, usually associated with trauma or injury with subsequent pain, whereas chronic LBP is often reported by patients in regular checkups and has led to a significant increase in the use of health care services over the past two decades. It is very important to differentiate between acute and chronic LBP in the clinical setting as these conditions—as well as their management and billing—are substantively different. Chronic back pain is generally treated with spinal injections [6,7], surgery [8,9], and/or pain medications [10,11], whereas anti-inflammatories and a rapid return to normal activities of daily living are generally the best recommendations for acute LBP [12].

However, acute and chronic LBP are usually not explicitly separated in electronic health records (EHRs) because of a lack of distinguishing codes. The International Classification of Diseases, 10th revision (ICD-10) standard only includes the code M54.5 to characterize *low back pain* diagnosis, and it does not provide modifiers to distinguish different LBP acuities [13]. Acuity is usually reported in clinical notes, requiring a retrospective chart review of the free text to characterize LBP events, which is time consuming and not scalable [14]. Moreover, acuity can be expressed in different ways. For example, the text could mention *acute low back pain* or *acute LBP*, but could also simply report *shooting pain down into the lower extremities*, *limited spine range of motion*, *vertebral tenderness*, *diffuse pain in lumbar muscles*, and so on [15]. This variability makes it difficult for clinical facilities and researchers to group LBP episodes by acuity to perform key tasks, such as defining appropriate diagnostic and billing codes; evaluating the effectiveness of prescribed treatments; and deriving data-driven therapeutic guidelines and improved diagnostic methods that could reduce time, disability, and cost.

This paper is the first to explore the use of automated approaches based on machine learning and information retrieval to analyze free-text clinical notes and identify the acuity of LBP episodes. Specifically, we use a set of manually annotated notes to train and evaluate various machine learning architectures based on logistic regression (LR), n-grams, topic models, word embeddings, and convolutional neural networks, and to demonstrate that some of these models are able to identify acute LBP episodes with promising precision. In addition, we demonstrate the ineffectiveness of using ICD-10 codes alone to train the models, reinforcing the idea that they are not sufficient to differentiate the acuity of LBP. Our overall objective was to build an automated framework that can help front line primary care providers (PCPs) in the development of targeted strategies and return-to-work (RTW) options for acute LBP episodes in clinical practice.

Background and Significance

PCPs are commonly the first medical practitioners to assess patient's musculoskeletal injuries and pain associated with these

injuries and are, therefore, in a unique position to offer reassurance, treatment options, and RTW recommendations catered to the acuity of the injury and pain associated with it. Several studies have documented increases in medication prescriptions and visits to physicians, physical therapists, and chiropractors for LBP episodes [16-18]. As individuals with chronic LBP seek care and use health care services more frequently than those with acute LBP, increases in health care use and costs for back pain are driven more by chronic than acute cases [19].

A rapid return to normal activities of daily living, including work, is generally the best activity recommendation for acute LBP management [12]. The number of workdays that are lost because of acute LBP can be reduced by implementing clinical practice guidelines in the primary care setting [20]. In previous work, Cruz et al built an RTW protocol tool for PCPs based on guidelines from the LBP literature [21]. On the basis of the type of work (eg, clerical, manual, or heavy) and the severity of the condition, the doctor would recommend RTW options (in partial or full duty capacity) within a certain number of days. The study found that physicians were likely to use this protocol, especially when it was integrated into the EHRs. However, the protocol was not always used for patients suffering from acute LBP as the research team was unable to quickly identify the acuity using only the structured EHR data (eg, ICD-10 codes). Acuity information was only available in the progress notes and was thus not incorporated into the automated recommendations. This prevented the research team from providing accurate feedback to PCPs based on a full picture of the patient's condition. A similar tool that could incorporate acuity information from notes could provide much more specific recommendations to PCPs that incorporate best practice guidelines for each acuity level. Besides leading to more precise care, this would streamline billing for LBP [22]. Similar needs arise for other musculoskeletal conditions, such as knee, elbow, and shoulder pain, where ICD-10 codes do not differentiate by pain level and acuity [23,24].

Machine learning methods for EHR data processing are enabling improved understanding of patient clinical trajectories, creating opportunities to derive new clinical insights [25,26]. In recent years, the application of deep learning, a hierarchical computational design based on layers of neural networks [27], to structured EHRs has led to promising results on clinical tasks such as disease phenotyping and prediction [28-33]. However, a wealth of relevant clinical information remains locked behind clinical narratives in the free text of notes. Natural language processing (NLP)—a branch of computer science that enables machines to process human language [34] for applications such as machine translation [35], text generation [36], and image captioning [37]—has been used to parse clinical notes to extract relevant insights that can guide clinical decisions [38]. Recent applications of deep learning to clinical NLP have classified clinical notes according to diagnosis or disease codes [39-41], predicted disease onset [32,42], and extracted primary cancer sites and their laterality in pathology reports [43,44]. However, although deep learning has successfully been applied to analyze clinical notes, traditional methods are still preferable when training data are limited [45,46].

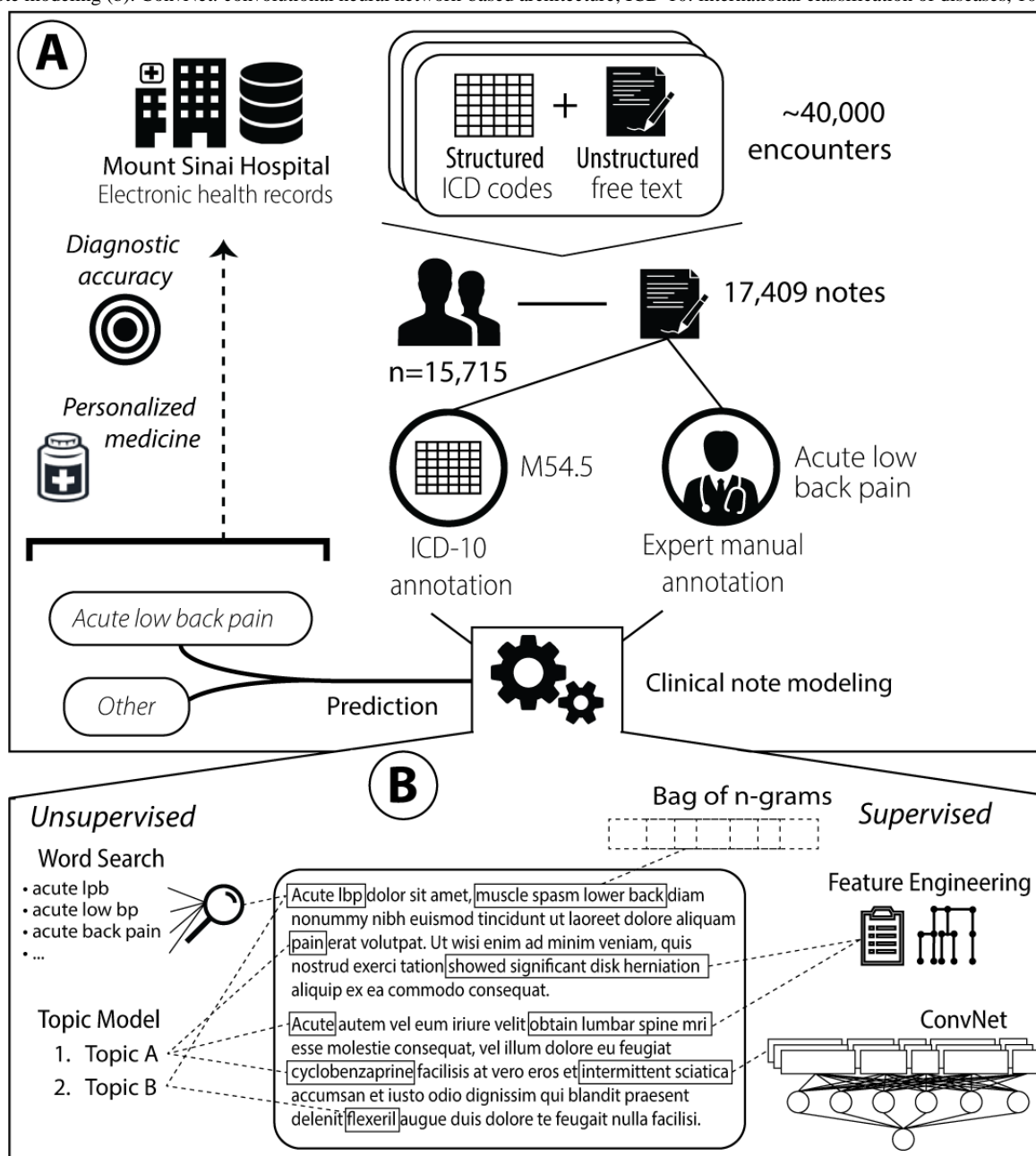
Regardless of the specific methodology, tools based on NLP applied to clinical narratives have not been widely used in clinical settings [31,38], despite the fact that physicians are likely to follow computer-assisted guidelines if recommendations are tied to their own observations [47]. In this paper, we present an NLP-based framework that can help physicians adhere to best practices and RTW recommendations for LBP. To the best of our knowledge, there are no studies to date that have applied machine learning to clinical notes to distinguish the acuity of a musculoskeletal condition in cases where it is not explicitly coded.

Methods

Overview

The conceptual steps of this study are summarized in Figure 1, specifically dataset composition, text processing, clinical notes modeling, and experimental evaluation. The overall goal was to evaluate the feasibility of automatically identifying clinical notes reporting acute LBP episodes.

Figure 1. Conceptual framework used to evaluate the use of automated approaches based on machine learning and information retrieval to analyze free-text clinical notes and identify acute low back pain episodes (a). The various unsupervised and supervised machine learning approaches used for clinical note modeling (b). ConvNet: convolutional neural network-based architecture; ICD-10: international classification of diseases, 10th revision.



Dataset

We used a set of free-text clinical notes extracted from the Mount Sinai data warehouse, made available for use under

institutional review board approval following Health Insurance Portability and Accountability Act guidelines. The Mount Sinai Health System is an urban tertiary care hospital located in the Upper East Side of Manhattan in New York City. It generates

a high volume of structured, semistructured, and unstructured data as part of its routine health care and clinical operations, which include inpatient, outpatient, and emergency room visits. These clinical notes were collected during a previous pilot study evaluating an RTW tool based on EHR data that included nearly 40,000 encounters for 15,715 patients spanning from 2016 to 2018 and clinical notes written by 81 different providers [21]. In that study, we used the published literature to develop a list of guidelines to determine the assessment and management of acute LBP episodes in clinical practice. In particular, we used ICD-10 codes and other parameters, such as *presenting complaint*, *pre-existing conditions*, *management factors*, and *imaging/radiology/test ordered*, to define and label the acuity of LBP in a clinical encounter. Following these guidelines, 14 individuals (physical medicine and rehabilitation fellows, residents, and medical students) manually reviewed a random set of 4291 clinical notes associated with these encounters and labeled all *acute low back pain* events. Each note was reviewed by at least two individuals and was further checked by a lead physician researcher if it was marked as ambiguous and/or there was discordance between reviewers.

This project leveraged the entire set of clinical notes that were collected in the previous study. In particular, we joined all the progress notes of these encounters under the same initial visit, and we eliminated duplicate, short (less than 3 words), and nonmeaningful reports. The final dataset was composed of 17,409 distinct clinical notes, with length ranging from 7 to 6638 words. Of this set, 3092 notes were manually reviewed in the previous study, and 891 of them were annotated as *acute LBP*. The remaining 14,317 notes were not manually evaluated and were related to different clinical domains, including various musculoskeletal disorders and potentially LBP events. In this final dataset, 1973 notes were also associated with an encounter billed with an ICD-10 M54.5 *Low back pain* code.

Text Processing

Every note in the dataset was tokenized, divided into sentences, and checked to remove punctuation; numbers; and nonrelevant concepts such as URLs, emails, and dates. Each note was then represented as a list of sentences, with every sentence being a list of lemmatized words represented as one-hot encodings. The vocabulary was composed of all the words appearing at least five times in the training set. The discarded words were corrected to the terms in the vocabulary having the minimum edit distance, that is, the minimum number of operations required to transform one string into the other [48]. This step reduced the number of misspelled words and prevented the accidental discarding of relevant information; at the same time, it also limited the size of the vocabulary to improve scalability [39]. Overall, the vocabulary covering the whole dataset comprised 56,142 unique words.

Clinical Note Modeling

We evaluated different approaches for identifying clinical notes that refer to acute LBP episodes. These included both supervised and unsupervised methods. Although we benefited from the use of high-quality manual annotations to train the supervised models, we also investigated alternatives that did not require manual annotation of notes. All these methods provided

straightforward explanations of their predictions, enabling us to validate each model and to identify parts of text and patterns that are relevant to the *acute LBP* predictions.

Keyword Search

We searched for a set of relevant keywords in the text. In particular, we looked for “acute low back pain,” “acute lbp,” “acute low bp,” and “acute back pain,” and we counted their occurrences in the text. We used the NegEx algorithm [49] to annotate and remove negated occurrences of the keywords. In the evaluation, we refer to this model as *WordSearch*.

Topic Modeling

We used topic modeling on the full set of words contained in the notes to capture abstract topics referred to in the dataset [50]. Topic modeling is an unsupervised inference process, in this case, implemented using latent Dirichlet allocation [51], which captures patterns of word co-occurrences within documents to define interpretable topics (ie, multinomial distribution of words) and represent a document as a multinomial over these topics. Every document can then be classified as talking about 1 or (usually) more topics. Topic modeling is often used in health care to generalize clinical notes, improve the automatic processing of patient data, and explore clinical datasets [52-55].

In this study, we assumed that 1 or more of these topics might refer to acute LBP. To discover them, we identified the most likely topics for a set of keywords (ie, “acute,” “low,” “back,” “pain,” “lbp,” and “bp”), and we manually reviewed them to retain only those that seemed more likely to characterize acute LBP episodes (ie, that included most of the keywords with high probability). We then considered the maximum likelihood among these topics as the probability that a report referred to acute LBP (ie, *TopicModel* in the experiments).

Bag of N-Grams

Each clinical note was represented as a bag of n-grams (BoN; with $n=1, \dots, 5$), with term frequency-inverse document frequency (TF-IDF) weights (determined from the corpus of documents). Each n-gram is a contiguous sequence of n words from the text. We considered all the words in the vocabulary and filtered the common stop words based on the English dictionary before building all the n-grams. The classification was implemented using LR with least absolute shrinkage and selection operator (LASSO; ie, *BoN-LR*).

Feature Engineering

We used the protocol built by Cruz et al [21] to define acute LBP episodes in the clinical notes. In particular, we used all the concepts described in that guideline, preprocessed them with the same algorithm used for the clinical notes, and built a set of 5154 distinct n-grams (with $n=1, \dots, 5$), which we refer to as *FeatEng*. We then represented each clinical note as a bag of *FeatEng* (ie, we counted the occurrences of only these n-grams in the text), normalized with TF-IDF weights, and classified them using LR with LASSO (ie, *FeatEng-LR*).

Deep Learning

We implemented an end-to-end deep neural network architecture based on convolutional neural networks that takes as input the full note and outputs its probability of being related to *acute LBP* (ie, *ConvNet* in the experiments). The first layer of the architecture maps the words to dense vector representations (ie, *embeddings*), which attempt to contextualize the semantic meaning of each word by creating a metric space where vectors of semantically similar words are close to each other. We applied word2vec with the skip-gram algorithm to the parsed notes [56] to initialize the embedding of each word in the vocabulary. Word2vec is commonly used with EHRs to learn embeddings of medical concepts from structured data and clinical notes [46,57-59].

The embeddings were then fed to a convolutional neural network inspired by the model described by Kim [60] and Liu et al [42]. This architecture concatenates representations of the text at different levels of abstraction by essentially choosing the most relevant n-grams at each level. Here, we first applied a set of parallel 1 dimensional (1D) convolutions on the input sequence with kernel sizes ranging from 1 to 5, thus simulating n-grams with $n=1, \dots, 5$. The outputs of each of these convolutions were then max-pooled over the whole sequence and concatenated to a $5 \times d$ dimensional vector, where d is the number of 1D convolutional filters. This representation was then fed to sequences of fully connected layers, which learn the interactions between the text features, and finally to a sigmoid layer that outputs the prediction probability.

The n-grams that are most relevant to the prediction, in this architecture, are those that activate the neurons in the max-pooling layer. Therefore, we used the log-odds that the n-gram contributes to the sigmoid decision function [42] as an indication of how much each n-gram influences the decision.

Evaluation Design

We evaluated all the architectures using a 10-fold cross-validation experiment, with every note appearing in the test set only once. In each training set, we used a random 90/10 split to train and validate all the model configurations. As baseline, we also report the results obtained by considering as *acute LBP* all the notes associated with the *Low back pain* M54.5 ICD-10 code (ie, *ICD-10* in the results).

Training Annotations

We considered 2 different sets of annotations as gold standards to train the supervised models. In the first experiment, we used the manually curated annotations provided with the dataset from previous work [21], whereas in the second experiment, we trained the models using the ICD-10 codes associated with each note encounter. Both experiments were evaluated using manual annotations. The rationale was to compare the feasibility of identifying acute LBP events when manual annotations are and are not available. We trained the classifier to output *acute LBP* versus *other* because the goal of the project was to identify clinical notes with acute LBP events rather than discriminate different facets of LBP events (eg, *chronic LBP* vs *acute LBP*).

Metrics

For all experiments, we report area under the receiver operating characteristic curve (AUC-ROC); precision, recall, and F score; and area under the precision-recall curve (AUC-PRC) [61]. The ROC curve is a plot of true positive rate versus false positive rate found over the set of predictions. F score is the harmonic mean of classification precision and recall per annotation, where precision is the number of correct positive results divided by the number of all positive results, and recall is the number of correct positive results divided by the number of positive results that should have been returned. The PRC is a plot of precision and recall for different thresholds. The areas under the ROC and PR curves are computed by integrating the corresponding curves.

Model Hyperparameters

The model hyperparameters were empirically tuned using the validation sets to optimize the results with both training annotations. In the topic modeling method, we inferred topics using the whole training set of documents and 200 topics (derived using perplexity analysis). Although seemingly more intuitive, using only the notes associated with the M54.5 *Low back pain* ICD-10 code actually produced worse results. For each fold, the most relevant topics associated with acute LBP were manually reviewed and used to annotate the notes. In the deep learning architecture, we used embeddings with size 300 and full-length notes. We trained word2vec just on the clinical note dataset to initialize embeddings. Preinitializing the embeddings with a general-purpose corpus did not lead to any improvement. Each convolutional neural network had 200 filters and used a rectified linear unit (ReLU) activation function. We added 2 fully connected layers of size 600 following the convolutional neural networks with ReLU activations and batch normalization. Dropout values across the layers were all set to 0.5. The architecture was trained using cross-entropy loss with the Adam optimizer for 5 epochs and batch size 32 (learning rate=0.001). The classification thresholds for precision, recall, and F score were found by ranging the value from 0.1 to 1, with 0.1 increments, and retaining, for each model, the value leading to the best results on the validation set.

Results

Table 1 and Figure 2 show the average results of the 10-fold cross-validation experiment for all the models considered. The best results were obtained by convolutional neural network-based architecture (ConvNet) when trained with the manual annotations. Although this is not entirely surprising given the success of deep learning for NLP when high-quality annotations and a large amount of data (ie, on the order of millions of training examples) are available, this was not certain in this domain where the training dataset was much smaller. As expected, the results obtained by the baseline and by training the models using the ICD-10 codes were not as good, confirming that the M54.5 ICD-10 code is not a sufficient indicator of acute LBP. TopicModel leads to similar performance but provides a more intuitive and potentially effective way for exploring the collection, extracting meaningful patterns that are related to acute LBP episodes. The most relevant topics included words

defining acute LBP (eg, acute, low, back, pain, lbp, spasm, lifting, sciatica) and also included several medications that are usually prescribed to treat inflammation and pain (eg, Cyclobenzaprine, Flexeril, and Advil). Although this approach might not be robust enough for clinical application, a refined and manually curated version of TopicModel promises to allow

an efficient prefiltering of clinical reports that can speed up the manual work required to annotate them. On the contrary, but as expected, WordSearch performed poorly as the condition is mentioned in too many different ways across the text, and simple keywords were not sufficient.

Table 1. The classification results in identifying clinical notes with acute low back pain (LBP) episodes averaged over the 10-fold cross-validation experiment. We compared different supervised and unsupervised strategies: keyword search (WordSearch), topic modeling (TopicModel), logistic regression with bag of n-grams (BoN-LR) and manual features (FeatEng-LR), and deep learning (ConvNet). The supervised models (ie, BoN-LR, FeatEng-LR, and ConvNet) were trained using manual annotations or M54.5 International Classification of Diseases, 10th revision (ICD-10) codes. The ICD-10 baseline simply considered as acute LBP all the notes associated with the generic M54.5 Low back pain ICD-10 code.

Model	Precision	Recall	F score	Area under the receiver operating characteristic curve	Area under the precision-recall curve
Baseline					
ICD-10 ^a	0.32	0.68	0.41	0.81	0.42
Unsupervised methods					
WordSearch	0.71	0.03	0.06	0.52	0.40
TopicModel	0.44	0.58	0.50	0.92	0.46
Trained with the M54.5 ICD-10 code					
BoN-LR ^b	0.50	0.70	0.59	0.83	0.42
FeatEng-LR ^c	0.47	0.59	0.52	0.88	0.41
ConvNet ^d	0.55	0.68	0.61	0.89	0.46
Trained with manual annotations					
BoN-LR	0.53	0.64	0.58	0.93	0.56
FeatEng-LR	0.58	0.66	0.62	0.93	0.58
ConvNet	0.65	0.73	0.70	0.98	0.72

^aICD-10: International Classification of Diseases, 10th revision codes.

^bBoN-LR: logistic regression with bag of n-grams.

^cFeatEng-LR: logistic regression with feature engineering.

^dConvNet: convolutional neural network-based architecture.

Figure 2. Receiver operating characteristic and precision-recall curves obtained when using as training data for BoN-LR, FeatEng-LR and ConvNet the manual annotations (a) and the M54.5 ICD-10 codes (b). ConvNet trained using the manual annotations obtained the best results. In the absence of manual annotations to use for training, TopicModel worked better than methods trained using ICD-10 codes, which proved not to be a good indicator to identify acuity in low back pain episodes. BoN-LR: logistic regression with bag of n-grams; ConvNet: convolutional neural network-based architecture; FeatEng-LR: logistic regression with feature engineering; ICD-10: international classification of diseases, 10th revision; PR: precision-recall; ROC: receiver operating characteristic.

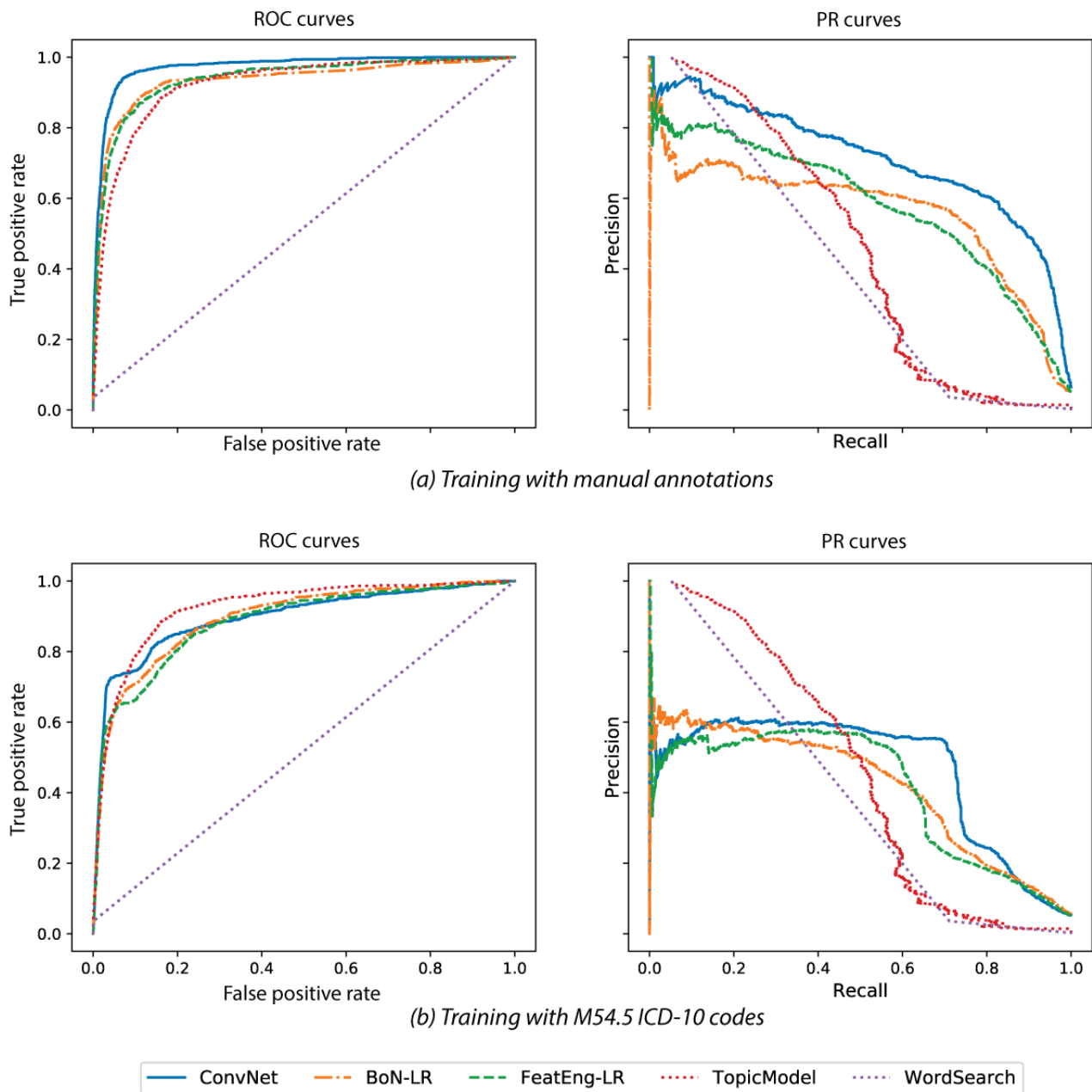


Figure 3 shows the classification results in terms of AUC-ROC and AUC-PRC when randomly subsampling the *acute LBP* manual annotations in the training set. We found that ConvNet always outperforms the other methods based on LR as well as TopicModel. In addition, we notice that using just 240 out of 800 (30.0%) manual annotations in the training set already leads

to better results than using ICD-10 codes as training labels. This is a particularly interesting insight as it shows that only minimal manual work is required to achieve good classifications; these can then be further improved by adding automatically annotated notes to the model (after manual verification) and retraining.

Figure 3. Area under the receiver operating characteristic and precision-recall curves obtained when training the supervised models using random subsamples of the manual annotations. TopicModel is reported as reference baseline. ConvNet obtained satisfactory results when trained using less manually annotated documents, showing robustness and scalability to the gold standard. AUC-PRC: area under the precision-recall curve; AUC-ROC: area under the receiver operating characteristic curve; BoN-LR: logistic regression with bag of n-grams; ConvNet: convolutional neural network-based architecture; FeatEng-LR: logistic regression with feature engineering.

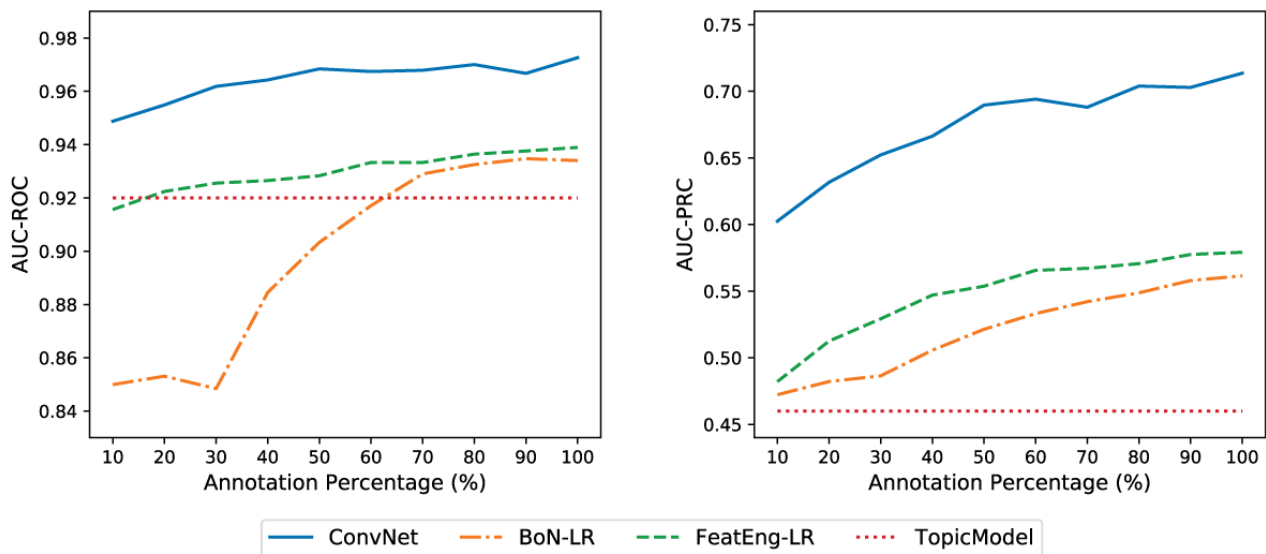
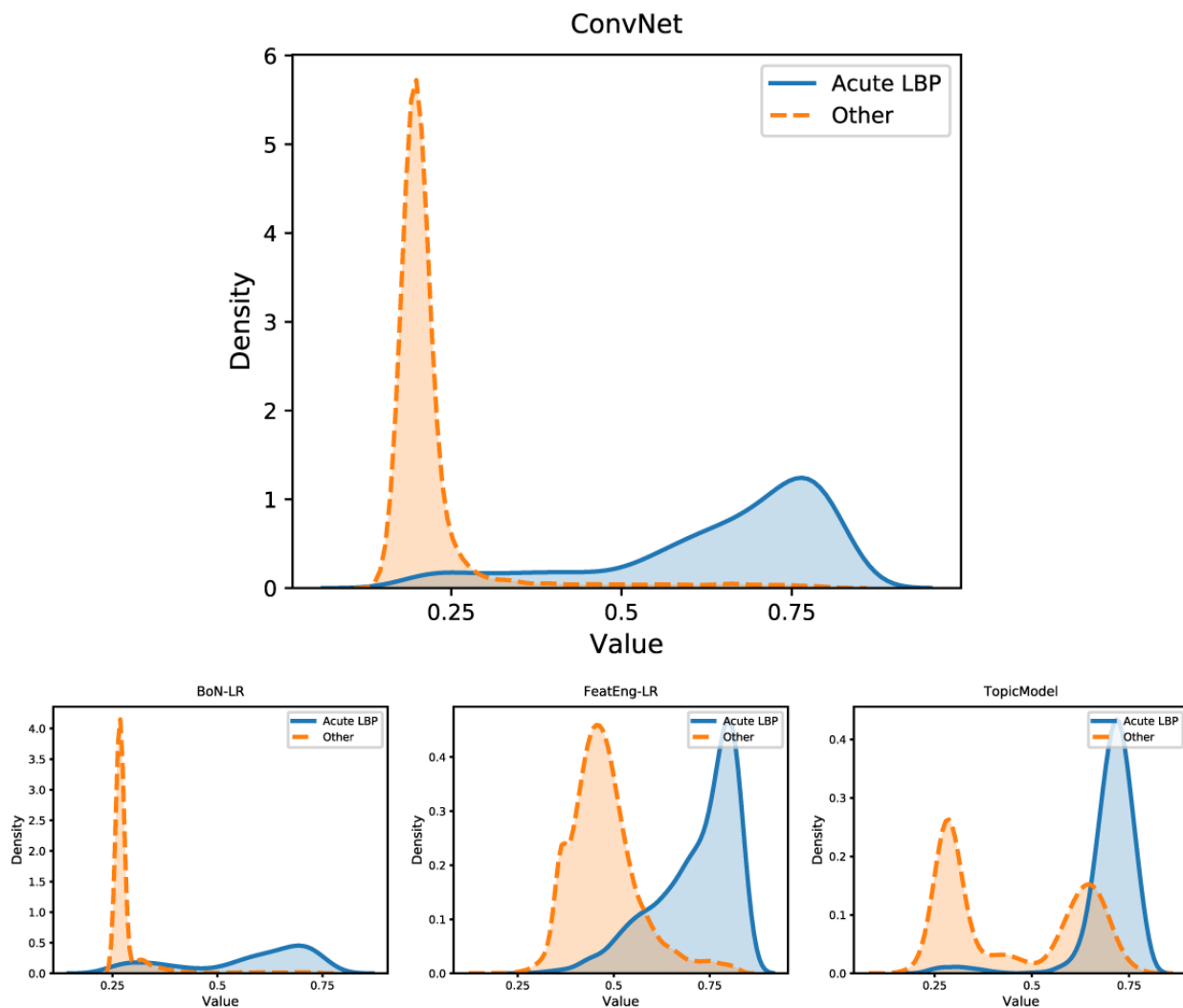


Figure 4 highlights the distributions of the classification scores (predicted probability of the label *acute LBP*) derived by several supervised models (trained with manual annotations) and TopicModel. ConvNet shows a clear separation between acute LBP notes and the rest of the dataset. In particular, all acute LBP notes had scores greater than 0.2, with 81.6% (727/891)

of them having scores greater than 0.5. On the contrary, only 347 controls had scores greater than 0.5, meaning that only a few notes were highly likely to be misclassified. Similarly, TopicModel had no controls with scores greater than 0.7, and all acute LBP notes had scores greater than 0.2.

Figure 4. Representation of the probability distribution of the scores obtained by BoN-LR, FeatEng-LR, ConvNet, and TopicModel. ConvNet led to a good separation between acute low back pain clinical notes and all the other documents. In other cases, such separation is not as clear, explaining the worse classification results obtained by those models. BoN-LR: logistic regression with bag of n-grams; ConvNet: convolutional neural network-based architecture; FeatEng-LR: logistic regression with feature engineering.



Finally, [Table 2](#) summarizes some of the n-grams driving the *acute LBP* predictions obtained by ConvNet (trained with manual annotations) across the experiments. Although some of these are obvious and refer to the disease itself (eg, “acute lbp”), others refer to medications (eg, “prescribed muscle relaxant”

and “flexeril”) and recommendations (eg, “rtw full duty quick”). Given their clinical meaning and relevance, all these patterns can be further analyzed and reviewed to potentially drive the development of guidelines for, for example, treatment and RTW options.

Table 2. Examples of n-grams that were relevant in identifying acute low back pain notes when using convolutional neural network-based architecture trained with manual annotations. The n-grams' relevance was determined by analyzing the neurons of the convolutional neural networks activating the max-pooling layers and their log-odds to contribute to the final output. Log-odds were filtered per notes and then averaged over all the notes and evaluation folds.

Type	Acute LBP ^a -related predictive n-grams
Diagnosis	<ul style="list-style-type: none"> • Muscle spasm lower back • Acute LBP flare • Been having acute back pain • Acute midline LBP • Sports acute bilateral LBP • Acute low back pain • Acute LBP
Related conditions	<ul style="list-style-type: none"> • Gait abnormality • Showed significant disk herniation • Intermittent sciatica • Spinal stenosis
Medications	<ul style="list-style-type: none"> • Back pain flare prescribed flexeril • Cyclobenzaprine • Flexeril • Naproxen for acute low back • Prescribed muscle relaxant
Recommendations	<ul style="list-style-type: none"> • Back brace for back pain • Obtain lumbar spine MRI^b • Recommendation RTW^c visit • RTW full duty quick

^aLBP: low back pain.

^bMRI: magnetic resonance imaging.

^cRTW: return-to-work.

Discussion

Principal Findings

In this work, we evaluated the use of several machine learning approaches to identify acute LBP episodes in free-text clinical notes to better personalize the treatment and management of this condition in primary care. The experimental results showed that it is possible to extract acute LBP episodes with promising precision, especially when at least some manually curated annotations are available. In this scenario, ConvNet, a deep learning architecture based on convolutional neural networks, significantly outperformed other shallow techniques based on BoN and LR, opening the possibility to boost performances using more complex architectures from current research in the NLP community. The implemented deep architecture also provides an easy mechanism to explain the predictions, leading to informed decision support based on model transparency [62,63] and the identification of meaningful patterns that can drive clinical decision making. If no annotations are available, experiments showed that the use of topic modeling is preferred to training a classifier using only the M54.5 ICD-10 codes (ie, *Low back pain*) associated with the clinical note encounter, which proved to be a poor indicator to discriminate LBP episodes. In addition, the topics identified can serve as an intuitive tool to inform guidelines and recommendations, to prefilter the documents, and to reduce the manual work required to annotate the notes. The proposed framework is inherently

domain agnostic and does not require any manual supervision to identify relevant features from the free text. Therefore, it can be leveraged in other musculoskeletal condition domains where acuity is not expressed in the ICD-10 diagnostic codes, such as knee, elbow, and shoulder pain.

Potential Applications

Medical care decisions are often based on heuristics and manually derived rule-based models constructed on previous knowledge and expertise [64]. Cognitive biases and personality traits, such as aversion to risk or ambiguity, overconfidence, and the anchoring effect, may lead to diagnostic inaccuracies and medical errors, resulting in mismanagement or inadequate utilization of resources [65]. In the LBP domain, this may lead to delays in finding the right therapy and assisting patients in the return to normal activities, increased risk of transitioning the condition from acute to chronic, discomfort for patients, and increased economic burdens on clinical facilities to adequately treat and manage this patient population. Deriving data-driven guidelines for treatment recommendations can help in reducing these cognitive biases and personality traits, leading to more consistent and accurate decisions. In this scenario, the proposed frameworks integrate seamlessly with the RTW tool proposed by Cruz et al [21] by including acuity-relevant information in the clinical notes and addressing 1 of the limitations of that study (ie, recommending the RTW tool at the point of care by accurately identifying the condition as acute LBP). Similarly, an understanding of the patterns driving the

predictions can lead to the development of new and improved treatment strategies for various types of injuries, which can be presented to the clinicians at the time of patient encounter to help them with better management of the condition. Although physicians will continue to have autonomy in determining optimal care pathways for their patients, the recommendations provided by the supporting framework will be useful to systematize and support their activities within the realm of the busy clinical practice. Posterior analysis of the clinical notes to infer acute LBP episodes can also help in assigning the proper diagnostic and billing codes for the encounter. In a foreseeable future scenario where, clinical observations are automatically transcribed via voice and EHRs are processed in real time, an automated tool that identifies acuity information could also improve the accuracy of diagnosis and billing in real time, with no need to wait for posterior evaluations.

Limitations

This work evaluated the feasibility of using machine learning to identify acute LBP episodes in clinical notes. Therefore, we compared different types of models (shallow vs deep) and learning frameworks (unsupervised vs supervised) to identify the best directions for implementation and deployment in real clinical settings. Although several of the architectures evaluated in this work obtained promising results, more sophisticated models are likely to improve these performances, especially in the deep learning domain. For example, algorithms based on attention models [66], Bidirectional Encoder Representations from Transformers [67], or XLNet [68] have shown encouraging results on similar NLP tasks and are likely to obtain better results in this domain as well. In this work, we only focused on processing clinical notes; however, embedding structured EHR data, especially medications, imaging studies, and/or laboratory tests, into the method should improve the results.

The dataset of clinical notes used in this study originated from a geographically diverse set of primary care clinics serving the New York City population across the city's metro area over a

limited period (ie, 2016 to 2018). Providers were enrolled and randomized into the study on a rolling basis, with the number of encounters for LBP varying for each individual provider, based on his/her own practice. The majority of the PCPs were assistant professors serving on the front lines. No specialists were included in the initial study, as the pilot project was only geared toward the PCPs. Consequently, the results of this study might not be applicable to specialty care practice.

Future Work

The classification of LBP episodes as acute or chronic at the point of care level within primary care practice is imperative for an RTW tool to be effectively used to render evidence-based guidelines. At this time, we plan to classify a large set of notes, derive patterns related to acute LBP, and extend the tool proposed by Cruz et al [21] according to them. We further plan to identify cases where the RTW tool can be easily deployed based on EHR integration in the clinical domain. We will also begin to address some of the methodological limitations of this study to optimize performance and evaluate its generalizability outside primary care. Finally, we aim to evaluate the feasibility of this type of approach for other musculoskeletal conditions, in particular, shoulder and knee pain.

Conclusions

This study demonstrates the feasibility of using machine learning to automatically identify acute LBP episodes from clinical reports using only unstructured free-text data. In particular, manually annotating a set of notes to use as a gold standard can lead to effective results, especially when using deep learning. Topic modeling can help in speeding up the annotation process, initiating an iterative process where initial predictions are validated and then used to refine and optimize the model. This approach provides a generalizable framework for learning to differentiate disease acuity in primary care, which can more accurately and specifically guide the diagnosis and treatment of LBP. It also provides a clear path toward improving the accuracy of coding and billing of clinical encounters for LBP.

Acknowledgments

IN and LC would like to thank the Pilot Projects Research Training Program of the New York and New Jersey Education and Research Center, National Institute for Occupational Safety and Health, for their funding (grant #T42 OH 008422). RM is grateful for the support from the Hasso Plattner Foundation and a courtesy GPU donation from NVIDIA.

Authors' Contributions

RM and IN initiated the idea and wrote the manuscript. IN collected the data and provided clinical support. RM conducted the research and the experimental evaluation. BP advised on evaluation strategies and refined the manuscript. BG, HL, and LC refined the manuscript. JD supported the research. All the authors edited and reviewed the manuscript.

Conflicts of Interest

None declared.

References

- Centers for Disease Control Prevention (CDC). Prevalence and most common causes of disability among adults--United States, 2005. *MMWR Morb Mortal Wkly Rep* 2009 May 1;58(16):421-426 [FREE Full text] [Medline: [19407734](#)]
- Ricci JA, Stewart WF, Chee E, Leotta C, Foley K, Hochberg MC. Back pain exacerbations and lost productive time costs in United States workers. *Spine (Phila Pa 1976)* 2006 Dec 15;31(26):3052-3060. [doi: [10.1097/01.brs.0000249521.61813.aa](#)] [Medline: [17173003](#)]

3. Luckhaupt SE, Dahlhamer JM, Gonzales GT, Lu ML, Groenewold M, Sweeney MH, et al. Prevalence, recognition of work-relatedness, and effect on work of low back pain among US workers. *Ann Intern Med* 2019 May 14. [doi: [10.7326/M18-3602](https://doi.org/10.7326/M18-3602)] [Medline: [31083729](https://pubmed.ncbi.nlm.nih.gov/31083729/)]
4. BMUS: The Burden of Musculoskeletal Diseases in the United States. Health Care Utilization and Economic Cost URL: <https://www.boneandjointburden.org/2014-report/if0/health-care-utilization-and-economic-cost> [accessed 2019-04-22]
5. Fairbank J, Gwilym SE, France JC, Daffner SD, Dettori J, Hermsmeyer J, et al. The role of classification of chronic low back pain. *Spine (Phila Pa 1976)* 2011 Oct 1;36(21 Suppl):S19-S42. [doi: [10.1097/BRS.0b013e31822ef72c](https://doi.org/10.1097/BRS.0b013e31822ef72c)] [Medline: [21952188](https://pubmed.ncbi.nlm.nih.gov/21952188/)]
6. Weiner DK, Kim YS, Bonino P, Wang T. Low back pain in older adults: are we utilizing healthcare resources wisely? *Pain Med* 2006;7(2):143-150. [doi: [10.1111/j.1526-4637.2006.00112.x](https://doi.org/10.1111/j.1526-4637.2006.00112.x)] [Medline: [16634727](https://pubmed.ncbi.nlm.nih.gov/16634727/)]
7. Friedly J, Chan L, Deyo R. Increases in lumbosacral injections in the Medicare population: 1994 to 2001. *Spine (Phila Pa 1976)* 2007 Jul 15;32(16):1754-1760. [doi: [10.1097/BRS.0b013e3180b9f96e](https://doi.org/10.1097/BRS.0b013e3180b9f96e)] [Medline: [17632396](https://pubmed.ncbi.nlm.nih.gov/17632396/)]
8. Deyo RA, Mirza SK. Trends and variations in the use of spine surgery. *Clin Orthop Relat Res* 2006 Feb;443:139-146. [doi: [10.1097/01.blo.0000198726.62514.75](https://doi.org/10.1097/01.blo.0000198726.62514.75)] [Medline: [16462438](https://pubmed.ncbi.nlm.nih.gov/16462438/)]
9. Deyo RA, Nachemson A, Mirza SK. Spinal-fusion surgery - the case for restraint. *N Engl J Med* 2004 Feb 12;350(7):722-726. [doi: [10.1056/NEJMs031771](https://doi.org/10.1056/NEJMs031771)] [Medline: [14960750](https://pubmed.ncbi.nlm.nih.gov/14960750/)]
10. Ballantyne JC. Opioids for the treatment of chronic pain: mistakes made, lessons learned, and future directions. *Anesth Analg* 2017 Nov;125(5):1769-1778. [doi: [10.1213/ANE.0000000000002500](https://doi.org/10.1213/ANE.0000000000002500)] [Medline: [29049121](https://pubmed.ncbi.nlm.nih.gov/29049121/)]
11. Luo X, Pietrobon R, Hey L. Patterns and trends in opioid use among individuals with back pain in the United States. *Spine (Phila Pa 1976)* 2004 Apr 15;29(8):884-90; discussion 891. [doi: [10.1097/00007632-200404150-00012](https://doi.org/10.1097/00007632-200404150-00012)] [Medline: [15082989](https://pubmed.ncbi.nlm.nih.gov/15082989/)]
12. Malmivaara A, Häkkinen U, Aro T, Heinrichs ML, Koskeniemi L, Kuosma E, et al. The treatment of acute low back pain--bed rest, exercises, or ordinary activity? *N Engl J Med* 1995 Feb 9;332(6):351-355. [doi: [10.1056/NEJM199502093320602](https://doi.org/10.1056/NEJM199502093320602)] [Medline: [7823996](https://pubmed.ncbi.nlm.nih.gov/7823996/)]
13. ICD-10 Data. 2020 ICD-10-CM Diagnosis Code M54.5: Low back pain URL: <https://www.icd10data.com/ICD10CM/Codes/M00-M99/M50-M54/M54-/M54.5> [accessed 2020-01-07]
14. Petersen T, Laslett M, Juhl C. Clinical classification in low back pain: best-evidence diagnostic rules based on systematic reviews. *BMC Musculoskelet Disord* 2017 May 12;18(1):188 [FREE Full text] [doi: [10.1186/s12891-017-1549-6](https://doi.org/10.1186/s12891-017-1549-6)] [Medline: [28499364](https://pubmed.ncbi.nlm.nih.gov/28499364/)]
15. Casazza BA. Diagnosis and treatment of acute low back pain. *Am Fam Physician* 2012 Feb 15;85(4):343-350 [FREE Full text] [Medline: [22335313](https://pubmed.ncbi.nlm.nih.gov/22335313/)]
16. Feuerstein M, Marcus SC, Huang GD. National trends in nonoperative care for nonspecific back pain. *Spine J* 2004;4(1):56-63. [doi: [10.1016/j.spinee.2003.08.003](https://doi.org/10.1016/j.spinee.2003.08.003)] [Medline: [14749194](https://pubmed.ncbi.nlm.nih.gov/14749194/)]
17. Kessler RC, Davis RB, Foster DF, van Rompay MI, Walters EE, Wilkey SA, et al. Long-term trends in the use of complementary and alternative medical therapies in the United States. *Ann Intern Med* 2001 Aug 21;135(4):262-268. [doi: [10.7326/0003-4819-135-4-200108210-00011](https://doi.org/10.7326/0003-4819-135-4-200108210-00011)] [Medline: [11511141](https://pubmed.ncbi.nlm.nih.gov/11511141/)]
18. Martin BI, Deyo RA, Mirza SK, Turner JA, Comstock BA, Hollingworth W, et al. Expenditures and health status among adults with back and neck problems. *J Am Med Assoc* 2008 Feb 13;299(6):656-664. [doi: [10.1001/jama.299.6.656](https://doi.org/10.1001/jama.299.6.656)] [Medline: [18270354](https://pubmed.ncbi.nlm.nih.gov/18270354/)]
19. Freburger JK, Holmes GM, Agans RP, Jackman AM, Darter JD, Wallace AS, et al. The rising prevalence of chronic low back pain. *Arch Intern Med* 2009 Feb 9;169(3):251-258 [FREE Full text] [doi: [10.1001/archinternmed.2008.543](https://doi.org/10.1001/archinternmed.2008.543)] [Medline: [19204216](https://pubmed.ncbi.nlm.nih.gov/19204216/)]
20. Rossignol M, Abenheim L, Séguin P, Neveu A, Collet JP, Ducruet T, et al. Coordination of primary health care for back pain. A randomized controlled trial. *Spine (Phila Pa 1976)* 2000 Jan 15;25(2):251-8; discussion 258. [doi: [10.1097/00007632-200001150-00018](https://doi.org/10.1097/00007632-200001150-00018)] [Medline: [10685491](https://pubmed.ncbi.nlm.nih.gov/10685491/)]
21. Cruz LC, Alamgir HA, Sheth P, Nabeel I. Development of a return to work tool for primary care providers for patients with low back pain: A pilot study. *J Family Med Prim Care* 2018;7(6):1185-1192 [FREE Full text] [doi: [10.4103/jfmpc.jfmpc_262_18](https://doi.org/10.4103/jfmpc.jfmpc_262_18)] [Medline: [30613495](https://pubmed.ncbi.nlm.nih.gov/30613495/)]
22. Owens JD, Hegmann KT, Thiese MS, Phillips AL. Impacts of adherence to evidence-based medicine guidelines for the management of acute low back pain on costs of worker's compensation claims. *J Occup Environ Med* 2019 Jun;61(6):445-452. [doi: [10.1097/JOM.0000000000001593](https://doi.org/10.1097/JOM.0000000000001593)] [Medline: [31167221](https://pubmed.ncbi.nlm.nih.gov/31167221/)]
23. Coding Strategies. Reporting Pain in ICD-10-CM URL: <https://www.codingstrategies.com/news/reporting-pain-icd-10-cm> [accessed 2020-01-07]
24. Gross DP, Armijo-Olivo S, Shaw WS, Williams-Whitt K, Shaw NT, Hartvigsen J, et al. Clinical decision support tools for selecting interventions for patients with disabling musculoskeletal disorders: a scoping review. *J Occup Rehabil* 2016 Sep;26(3):286-318 [FREE Full text] [doi: [10.1007/s10926-015-9614-1](https://doi.org/10.1007/s10926-015-9614-1)] [Medline: [26667939](https://pubmed.ncbi.nlm.nih.gov/26667939/)]
25. Jensen PB, Jensen LJ, Brunak S. Mining electronic health records: towards better research applications and clinical care. *Nat Rev Genet* 2012 May 2;13(6):395-405. [doi: [10.1038/nrg3208](https://doi.org/10.1038/nrg3208)] [Medline: [22549152](https://pubmed.ncbi.nlm.nih.gov/22549152/)]

26. Glicksberg BS, Johnson KW, Dudley JT. The next generation of precision medicine: observational studies, electronic health records, biobanks and continuous monitoring. *Hum Mol Genet* 2018 May 1;27(R1):R56-R62. [doi: [10.1093/hmg/ddy114](https://doi.org/10.1093/hmg/ddy114)] [Medline: [29659828](https://pubmed.ncbi.nlm.nih.gov/29659828/)]
27. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature* 2015 May 28;521(7553):436-444. [doi: [10.1038/nature14539](https://doi.org/10.1038/nature14539)] [Medline: [26017442](https://pubmed.ncbi.nlm.nih.gov/26017442/)]
28. Miotto R, Li L, Kidd BA, Dudley JT. Deep Patient: an unsupervised representation to predict the future of patients from the electronic health records. *Sci Rep* 2016 May 17;6:26094 [FREE Full text] [doi: [10.1038/srep26094](https://doi.org/10.1038/srep26094)] [Medline: [27185194](https://pubmed.ncbi.nlm.nih.gov/27185194/)]
29. Miotto R, Wang F, Wang S, Jiang X, Dudley JT. Deep learning for healthcare: review, opportunities and challenges. *Brief Bioinform* 2018 Nov 27;19(6):1236-1246 [FREE Full text] [doi: [10.1093/bib/bbx044](https://doi.org/10.1093/bib/bbx044)] [Medline: [28481991](https://pubmed.ncbi.nlm.nih.gov/28481991/)]
30. Choi E, Bahadori MT, Schuetz A, Stewart WF, Sun J. Doctor AI: predicting clinical events via recurrent neural networks. *JMLR Workshop Conf Proc* 2016 Aug;56:301-318 [FREE Full text] [Medline: [28286600](https://pubmed.ncbi.nlm.nih.gov/28286600/)]
31. Xiao C, Choi E, Sun J. Opportunities and challenges in developing deep learning models using electronic health records data: a systematic review. *J Am Med Inform Assoc* 2018 Oct 1;25(10):1419-1428 [FREE Full text] [doi: [10.1093/jamia/ocy068](https://doi.org/10.1093/jamia/ocy068)] [Medline: [29893864](https://pubmed.ncbi.nlm.nih.gov/29893864/)]
32. Rajkumar A, Oren E, Chen K, Dai AM, Hajaj N, Hardt M, et al. Scalable and accurate deep learning with electronic health records. *NPJ Digit Med* 2018;1:18 [FREE Full text] [doi: [10.1038/s41746-018-0029-1](https://doi.org/10.1038/s41746-018-0029-1)] [Medline: [31304302](https://pubmed.ncbi.nlm.nih.gov/31304302/)]
33. Miotto R, Li L, Dudley JT. Deep learning to predict patient future diseases from the electronic health records. In: Ferro N, Crestani F, Moens MF, editors. *Advances in Information Retrieval*. Cham: Springer; 2016:768-774.
34. Goldberg Y. A primer on neural network models for Natural Language Processing. *J Artif Intell Res* 2016;57(7):345-420. [doi: [10.1613/jair.4992](https://doi.org/10.1613/jair.4992)]
35. Wu Y, Schuster M, Chen Z, Le QV, Norouzi M, Macherey W, et al. arXiv e-Print archive. 2016. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation URL: <http://arxiv.org/abs/1609.08144> [accessed 2020-01-07]
36. Kannan A, Kurach K, Ravi S, Kaufmann T, Tomkins A, Miklos B, et al. Smart Reply: Automated Response Suggestion for Email. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. San Francisco, CA: ACM; 2016 Presented at: KDD'16; August 13-17, 2016; San Francisco, United States p. 955-964.
37. Vinyals O, Toshev A, Bengio S, Erhan D. Show and Tell: A Neural Image Caption Generator. In: *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition*. 2015 Presented at: CVPR'15; June 7-12, 2015; Boston, MA URL: https://www.cv-foundation.org/openaccess/content_cvpr_2015/papers/Vinyals_Show_and_Tell_2015_CVPR_paper.pdf
38. Sheikhalishahi S, Miotto R, Dudley JT, Lavelli A, Rinaldi F, Osmani V. Natural Language Processing of clinical notes on chronic diseases: systematic review. *JMIR Med Inform* 2019 Apr 27;7(2):e12239 [FREE Full text] [doi: [10.2196/12239](https://doi.org/10.2196/12239)] [Medline: [31066697](https://pubmed.ncbi.nlm.nih.gov/31066697/)]
39. Baumel T, Nassour-Kassis J, Cohen R, Elhadad M, Elhadad N. arXiv e-Print archive. 2017. Multi-Label Classification of Patient Notes a Case Study on ICD Code Assignment URL: <http://arxiv.org/abs/1709.09587> [accessed 2020-01-07]
40. Mullenbach J, Wiegrefe S, Duke J, Sun J, Eisenstein J. arXiv e-Print archive. 2018. Explainable Prediction of Medical Codes from Clinical Text URL: <http://arxiv.org/abs/1802.05695> [accessed 2020-01-07]
41. Shi H, Xie P, Hu Z, Zhang M, Xing EP. arXiv e-Print archive. 2017. Towards Automated ICD Coding Using Deep Learning URL: <http://arxiv.org/abs/1711.04075> [accessed 2020-01-07]
42. Liu J, Zhang Z, Razavian N. arXiv e-Print archive. 2018. Deep EHR: Chronic Disease Prediction Using Medical Notes URL: <http://arxiv.org/abs/1808.04928> [accessed 2020-01-07]
43. Yoon HJ, Ramanathan A, Tourassi G. Multi-task deep neural networks for automated extraction of primary site and laterality information from cancer pathology reports. In: Angelov P, Manolopoulos Y, Iliadis L, Roy A, Vellasco M, editors. *Advances in Big Data*. Cham, Switzerland: Springer; 2016:195-204.
44. Qiu JX, Yoon HJ, Fearn PA, Tourassi GD. Deep learning for automated extraction of primary sites from cancer pathology reports. *IEEE J Biomed Health Inform* 2018 Jan;22(1):244-251. [doi: [10.1109/JBHI.2017.2700722](https://doi.org/10.1109/JBHI.2017.2700722)] [Medline: [28475069](https://pubmed.ncbi.nlm.nih.gov/28475069/)]
45. Turner CA, Jacobs AD, Marques CK, Oates JC, Kamen DL, Anderson PE, et al. Word2Vec inversion and traditional text classifiers for phenotyping lupus. *BMC Med Inform Decis Mak* 2017 Aug 22;17(1):126 [FREE Full text] [doi: [10.1186/s12911-017-0518-1](https://doi.org/10.1186/s12911-017-0518-1)] [Medline: [28830409](https://pubmed.ncbi.nlm.nih.gov/28830409/)]
46. Gehrman S, Deroncourt F, Li Y, Carlson ET, Wu JT, Welt J, et al. arXiv e-Print archive. 2017. Comparing Rule-Based and Deep Learning Models for Patient Phenotyping URL: <http://arxiv.org/abs/1703.08705> [accessed 2020-01-07]
47. Davis DA, Taylor-Vaisey A. Translating guidelines into practice. A systematic review of theoretic concepts, practical experience and research evidence in the adoption of clinical practice guidelines. *Can Med Assoc J* 1997 Aug 15;157(4):408-416 [FREE Full text] [Medline: [9275952](https://pubmed.ncbi.nlm.nih.gov/9275952/)]
48. Navarro G. A guided tour to approximate string matching. *ACM Comput Surv* 2001;33(1):31-88 [FREE Full text] [doi: [10.1145/375360.375365](https://doi.org/10.1145/375360.375365)]
49. Chapman WW, Bridewell W, Hanbury P, Cooper GF, Buchanan BG. A simple algorithm for identifying negated findings and diseases in discharge summaries. *J Biomed Inform* 2001 Oct;34(5):301-310 [FREE Full text] [doi: [10.1006/jbin.2001.1029](https://doi.org/10.1006/jbin.2001.1029)] [Medline: [12123149](https://pubmed.ncbi.nlm.nih.gov/12123149/)]

50. Blei DM. Probabilistic topic models. *Commun ACM* 2012;55(4):77. [doi: [10.1145/2133806.2133826](https://doi.org/10.1145/2133806.2133826)]
51. Blei DM, Ng AY, Jordan MI. Latent Dirichlet Allocation. *J Mach Learn Res* 2003;3:993-1022 [FREE Full text]
52. Miotto R, Weng C. Case-based reasoning using electronic health records efficiently identifies eligible patients for clinical trials. *J Am Med Inform Assoc* 2015 Apr;22(e1):e141-e150 [FREE Full text] [doi: [10.1093/jamia/ocu050](https://doi.org/10.1093/jamia/ocu050)] [Medline: [25769682](https://pubmed.ncbi.nlm.nih.gov/25769682/)]
53. Perotte AJ, Wood F, Elhadad N, Bartlett N. Hierarchically Supervised Latent Dirichlet Allocation. In: *Proceedings of the Advances in Neural Information Processing Systems 24*. New York, NY: Curran Associates, Inc; 2011 Presented at: NIPS'11; December 12-17, 2011; Granada, Spain p. 2609-2617 URL: <https://papers.nips.cc/paper/4313-hierarchically-supervised-latent-dirichlet-allocation.pdf> [doi: [10.1108/09504120310455975](https://doi.org/10.1108/09504120310455975)]
54. Chang KR, Lou X, Karaletsos T, Crosbie C, Gardos S, Artz D, et al. bioRxiv - the preprint server for Biology. 2016. An Empirical Analysis of Topic Modeling for Mining Cancer Clinical Notes URL: <https://www.biorxiv.org/content/10.1101/062307v1> [accessed 2020-01-07]
55. Cohen R, Aviram I, Elhadad M, Elhadad N. Redundancy-aware topic modeling for patient record notes. *PLoS One* 2014;9(2):e87555 [FREE Full text] [doi: [10.1371/journal.pone.0087555](https://doi.org/10.1371/journal.pone.0087555)] [Medline: [24551060](https://pubmed.ncbi.nlm.nih.gov/24551060/)]
56. Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J. Distributed Representations of Words and Phrases and their Compositionality. In: *Proceedings of the Advances in Neural Information Processing Systems 26*. New York, NY: Curran Associates, Inc; 2013 Presented at: NIPS'13; December 5-10, 2013; Lake Tahoe, Nevada, USA p. 3111-3119 URL: <https://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf>
57. Glicksberg BS, Miotto R, Johnson KW, Shameer K, Li L, Chen R, et al. Automated disease cohort selection using word embeddings from Electronic Health Records. *Pac Symp Biocomput* 2018;23:145-156 [FREE Full text] [Medline: [29218877](https://pubmed.ncbi.nlm.nih.gov/29218877/)]
58. Choi Y, Chiu CY, Sontag D. Learning low-dimensional representations of medical concepts. *AMIA Jt Summits Transl Sci Proc* 2016;2016:41-50 [FREE Full text] [Medline: [27570647](https://pubmed.ncbi.nlm.nih.gov/27570647/)]
59. Wang Y, Liu S, Afzal N, Rastegar-Mojarad M, Wang L, Shen F, et al. A comparison of word embeddings for the biomedical natural language processing. *J Biomed Inform* 2018 Nov;87:12-20 [FREE Full text] [doi: [10.1016/j.jbi.2018.09.008](https://doi.org/10.1016/j.jbi.2018.09.008)] [Medline: [30217670](https://pubmed.ncbi.nlm.nih.gov/30217670/)]
60. Kim Y. arXiv e-Print archive. 2014. Convolutional Neural Networks for Sentence Classification URL: <http://arxiv.org/abs/1408.5882> [accessed 2020-01-07]
61. Baeza-Yates R, Ribeiro-Neto B. *Modern Information Retrieval: The Concepts and Technology behind Search*. Second Edition. Boston, MA: Addison-Wesley Professional; 2011.
62. Holzinger A, Biemann C, Pattichis CS, Kell DB. arXiv e-Print archive. 2017. What Do We Need to Build Explainable AI Systems for the Medical Domain? URL: <http://arxiv.org/abs/1712.09923> [accessed 2020-01-07]
63. Lipton ZC. arXiv e-Print archive. 2016. The Mythos of Model Interpretability URL: <http://arxiv.org/abs/1606.03490> [accessed 2020-01-07]
64. Marewski JN, Gigerenzer G. Heuristic decision making in medicine. *Dialogues Clin Neurosci* 2012 Mar;14(1):77-89 [FREE Full text] [Medline: [22577307](https://pubmed.ncbi.nlm.nih.gov/22577307/)]
65. Saposnik G, Redelmeier D, Ruff CC, Tobler PN. Cognitive biases associated with medical decisions: a systematic review. *BMC Med Inform Decis Mak* 2016 Nov 3;16(1):138 [FREE Full text] [doi: [10.1186/s12911-016-0377-1](https://doi.org/10.1186/s12911-016-0377-1)] [Medline: [27809908](https://pubmed.ncbi.nlm.nih.gov/27809908/)]
66. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is All you Need. In: *Proceedings of the Advances in Neural Information Processing Systems 30*. New York, NY: Curran Associates, Inc; 2017 Presented at: NIPS'17; December 4-9, 2017; Long Beach, CA, USA p. 5998-6008 URL: <https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>
67. Devlin J, Chang MW, Lee K, Toutanova K. arXiv e-Print archive. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding URL: <http://arxiv.org/abs/1810.04805> [accessed 2020-01-07]
68. Yang Z, Dai Z, Yang Y, Carbonell J, Salakhutdinov R, Le QV. arXiv e-Print archive. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding URL: <http://arxiv.org/abs/1906.08237> [accessed 2020-01-07]

Abbreviations

- 1D:** 1 dimensional
- AUC-PRC:** area under the precision-recall curve
- AUC-ROC:** area under the receiver operating characteristic curve
- BoN:** bag of n-grams
- ConvNet:** convolutional neural network-based architecture
- EHR:** electronic health record
- ICD-10:** International Classification of Diseases, 10th revision
- LASSO:** least absolute shrinkage and selection operator
- LBP:** low back pain
- LR:** logistic regression

NLP: natural language processing
PCP: primary care provider
ReLu: rectified linear unit
RTW: return-to-work
TF-IDF: term frequency-inverse document frequency

Edited by C Lovis; submitted 01.11.19; peer-reviewed by V Osmani, M Boukhechba; accepted 15.12.19; published 27.02.20

Please cite as:

Miotto R, Percha BL, Glicksberg BS, Lee HC, Cruz L, Dudley JT, Nabeel I

Identifying Acute Low Back Pain Episodes in Primary Care Practice From Clinical Notes: Observational Study

JMIR Med Inform 2020;8(2):e16878

URL: <http://medinform.jmir.org/2020/2/e16878/>

doi: [10.2196/16878](https://doi.org/10.2196/16878)

PMID:

©Riccardo Miotto, Bethany L Percha, Benjamin S Glicksberg, Hao-Chih Lee, Lisanne Cruz, Joel T Dudley, Ismail Nabeel. Originally published in JMIR Medical Informatics (<http://medinform.jmir.org>), 27.02.2020. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in JMIR Medical Informatics, is properly cited. The complete bibliographic information, a link to the original publication on <http://medinform.jmir.org/>, as well as this copyright and license information must be included.