

Annual Review of Biomedical Data Science

Modern Clinical Text Mining: A Guide and Review

Bethany Percha

Department of Medicine and Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY 10025, USA; email: bethany.percha@mssm.edu

Annu. Rev. Biomed. Data Sci. 2021. 4:165–87

First published as a Review in Advance on
May 26, 2021

The *Annual Review of Biomedical Data Science* is
online at biomedataci.annualreviews.org

<https://doi.org/10.1146/annurev-biomedataci-030421-030931>

Copyright © 2021 by Annual Reviews.
All rights reserved

ANNUAL
REVIEWS **CONNECT**

www.annualreviews.org

- Download figures
- Navigate cited references
- Keyword search
- Explore related articles
- Share via email or social media

Keywords

text mining, natural language processing, electronic health record, clinical text, machine learning

Abstract

Electronic health records (EHRs) are becoming a vital source of data for healthcare quality improvement, research, and operations. However, much of the most valuable information contained in EHRs remains buried in unstructured text. The field of clinical text mining has advanced rapidly in recent years, transitioning from rule-based approaches to machine learning and, more recently, deep learning. With new methods come new challenges, however, especially for those new to the field. This review provides an overview of clinical text mining for those who are encountering it for the first time (e.g., physician researchers, operational analytics teams, machine learning scientists from other domains). While not a comprehensive survey, this review describes the state of the art, with a particular focus on new tasks and methods developed over the past few years. It also identifies key barriers between these remarkable technical advances and the practical realities of implementation in health systems and in industry.

Information

extraction: extracting structured information (e.g., concepts, relations, events) from unstructured text; examples include named-entity recognition, concept normalization, and relation extraction; often considered a subdomain of NLP

1. INTRODUCTION

Among the most significant barriers to large-scale deployment of electronic health records (EHRs) in quality improvement, operations, and research is the amount of EHR data stored as unstructured text (1). Structured, machine-computable data, such as procedure and diagnosis codes, are in the minority. The bulk of information relating clinical findings to decisions and communicating the logical and deductive processes of medicine is buried within progress notes, radiology and pathology reports, and other free-text documents (2, 3). Examples include treatment goals and outcomes (e.g., success or failure of treatments, criteria for success, decisions about subsequent treatments); interpretations of radiology and pathology images and laboratory test results; social determinants of health (e.g., social connection/isolation, housing issues, mentions of financial resource strain) (4); symptoms, symptom changes, and their interpretation (5); past medical history and family history; patients' emotional disposition, mood, and interactions with health providers; detailed descriptions of procedures (e.g., labor and delivery, heart catheterization, imaging studies, surgeries); adherence to treatment plans (e.g., medications, physical therapy, procedures); allergies, side effects, and other adverse events; results of physical examination (e.g., review of systems and interpretation of findings); patients' reasons for seeing a health provider; primary and secondary complaints; psychiatric evaluations and records of therapy sessions; and discharge summaries and follow-up plans.

Some have speculated that modern machine learning algorithms, combined with EHR and other patient data, will enable the convergence of human and machine intelligence in healthcare (6, 7). From a practical standpoint, such a vision hinges on text mining. Without the ability to reliably process and interpret vast quantities of clinical text, all attempts to create high-performance predictive models, phenotyping algorithms, and data-driven treatment strategies (precision medicine) will face substantial challenges.

For the past several decades, a community of researchers working at the intersection of computer science and medicine has developed strategies for information extraction and modeling of clinical text, using techniques somewhat distinct from those of the broader natural language processing (NLP) research community (8, 9). Their efforts have led to the development of new methods and the production of both commercial (10) and open-source (11) software systems for clinical text mining. In recent years, technology giants like Amazon and Google have also recognized the importance of clinical text mining and joined the fray; Amazon Comprehend Medical (12) now comes packaged as a software add-on to Amazon Web Services, incentivizing storage of EHR data on Amazon's HIPAA (Health Insurance Portability and Accountability Act)-compliant cloud platform by providing seamless clinical text processing. Dedicated clinical text-processing companies such as (as of this writing) Clinithink (<http://clinithink.com>), Linguamatics (<http://linguamatics.com>), and Apixio (<http://apixio.com>) have built proprietary systems of their own, promising to improve clinical trial recruitment, disease registry creation, government reporting, and billing, all through improved mining of unstructured clinical text.

As a data scientist with a background in biomedical text mining, I am frequently approached by physician colleagues and academic and industry collaborators who, for various reasons, have found themselves needing to process clinical text. Many perceive clinical text mining to be a solved problem and believe that one can simply apply a packaged clinical NLP system to extract structured data for a variety of downstream applications. As a result, I often find myself explaining the limits of current NLP technology and the fact that clinical NLP encompasses many different goals, progress on some of which is further along than others. The purpose of this review, therefore, is to provide a starting point for those who are encountering clinical text mining for the first time. Far from a comprehensive survey, it focuses on a subset of methods and ideas that are particularly

clear, generalizable, and useful as starting points for further explorations of the field. Importantly, nothing I discuss here requires access to institution-specific or proprietary software, rule sets, or training corpora. My goal is to provide outsiders with a realistic baseline for what it is possible to accomplish with clinical text mining today.

2. A SHORT TAXONOMY OF TASKS AND APPROACHES

2.1. Information Extraction Versus Modeling

Clinical text can play multiple roles in a project, so it is important to start by defining one's overall goal and how the text fits in. For example, electronic phenotyping algorithms (13–16) often combine clinical notes with structured data, such as diagnosis codes, medication orders, and procedures, to make a prediction about whether a patient has a disease or other phenotype. Here the primary goal of text mining is information extraction: converting the text into a set of structured features that can be combined with other types of features to produce an answer (17). EHR search indexing, knowledgebase construction, and patient timeline building are similar in their focus on information extraction.

Equally important, however, are problems where the goal is to make a prediction or inference from the text itself—for example, to classify mammography reports by BI-RADS (Breast Imaging Reporting and Data System) category (18) or to cluster clinical documents to uncover latent structure (19). This may or may not require a separate information extraction step. For example, methods such as end-to-end, deep learning–based text classification models (20–22), which produce answers directly from the raw text, often shine in such cases. One important consideration is whether human-interpretable features are necessary for the project or whether the algorithm can be allowed to learn its own representations of text automatically in the course of solving a downstream task (22).

2.2. Rule-Based Versus Statistical Approaches

Clinical NLP systems fall into two broad categories: rule-based and statistical. Rule-based NLP systems codify expert knowledge into a set of structured rules, or templates, that produce structured information when applied to unstructured text. For example, a rule might specify patterns of words, phrases, or parts of speech that signal the presence of a particular type of entity: “If the word ‘received’ is followed by a noun, followed by ‘for’ and then a disease name, assume the noun is a drug name.” Many of the best-performing clinical NLP systems are rule based: Of 263 clinical text mining articles reviewed by Wang et al. in 2018 (17), 171 (65%) used rule-based methods. However, rule-based systems have two important disadvantages. First, domain experts must often expend substantial time and effort to construct the rules. Second, because they are domain specific, they do not generalize well to new problems; a rule-based system for identifying drug names in text will not be good at anything other than identifying drug names in text.

The alternative is a statistical NLP system built using a statistical learning (i.e., machine learning) algorithm. If provided with some text in which all of the drug names are labeled, for example, the algorithm will try to identify patterns that indicate that a particular span of text is a drug name (9, Ch. 8). Learning algorithms themselves are often task independent, which is one of their key advantages. However, statistical learning algorithms require annotated training data, which in the clinical domain is often limited or nonexistent (23). In addition, privacy concerns often make it impossible to share training data across institutions. As a result, while the NLP community has increasingly turned toward machine learning and away from rules-based approaches, clinical text mining maintains a strong focus on rules (8).

Electronic

phenotyping:

identifying patients with characteristics of interest (e.g., exposures, diseases, outcomes), usually from EHRs, claims, or other administrative data; also called cohort identification

Text classification:

assigning a label, or category, to text based on its content; examples include document classification (e.g., of radiology, pathology, or autopsy reports) and sentence classification

Rule-based NLP

system: applies a set of expert-defined rules, or templates, to perform an NLP task; the downside is the need for expert time and effort

Statistical NLP

system: learns how to perform an NLP task by applying machine learning algorithms to training data; the downside is the need for (often large amounts of) training data

Named-entity recognition (NER): identifying and locating mentions of conceptual categories, such as drug, symptom, or disease names, in text

Concept normalization: assigning a unique identity to an entity name recognized in text; in the biomedical domain, names are typically mapped to concepts from structured terminologies or ontologies

3. SOFTWARE FOR CLINICAL INFORMATION EXTRACTION

The three most common information extraction tasks—named-entity recognition (NER) (37–39), concept normalization (40, 41), and relation extraction (Section 7)—are still active areas of research. However, because of the broad need for basic information extraction in applied tasks like medical coding, search, and case finding, software systems have been developed to perform these tasks automatically. This section reviews current state-of-the-art methods and systems and provides examples of the type of output one can expect from each system. A summary of software packages for clinical text mining can be found in **Table 1**.

3.1. Named-Entity Recognition

NER is the task of identifying and locating mentions of conceptual categories, such as drug, symptom, or disease names, in text. It is perhaps the most widely studied information extraction task, and existing clinical NER systems can already identify a variety of clinically relevant entities, including problems, tests, and treatments (32, 37, 38); medication and adverse event names (42, 43); and protected health information (44, 45). **Figure 1** shows the raw output from Stanza (26, 32), a state-of-the-art NER system trained to tag clinical entities using `test`, `problem`, and `treatment` concept annotations from the 2010 i2b2/VA (Informatics for Integrating Biology and the Bedside/ Veterans Affairs) dataset (46).

Table 1 Clinical text mining software and resources

Resource	Language	URL	Reference
NLTK Toolkit	Python	http://nltk.org	24
Stanford CoreNLP	Java	http://stanfordnlp.github.io/CoreNLP	25
Stanza	Python	http://stanfordnlp.github.io/stanza	26
spaCy	Python, Cython	http://spacy.io	
scispaCy	Python	http://allenai.github.io/scispacy	27
Apache OpenNLP	Java	http://opennlp.apache.org	
CRFSuite	Python	http://chokkan.org/software/crfsuite	
scikit-learn	Python	http://scikit-learn.org (text preprocessing: <code>sklearn.feature_extraction.text</code>)	
Gensim	Python	http://radimrehurek.com/gensim/index.html	28
BERT	Python	http://github.com/google-research/bert	29
MetaMap	Java	http://metamap.nlm.nih.gov	30
MetaMap Lite	Java	http://metamap.nlm.nih.gov/MetaMapLite.shtml	31
cTAKES	Java	http://ctakes.apache.org	11
Stanza clinical	Python	http://stanza.run/bio	32
DNorm	Java, REST API	http://ncbi.nlm.nih.gov/research/bionlp/Tools/dnorm	33
Clinical BERT	Python	http://github.com/EmilyAlsentzer/clinicalBERT	34
	Python	http://github.com/kexinhuang12345/clinicalBERT	35
UMLS	NA (extraction software in Java)	http://nlm.nih.gov/research/umls/index.html	36

The above tools are popular choices for general and clinical text processing (e.g., word and sentence tokenization, part-of-speech tagging, chunking, parsing, NER, word and phrase embeddings). The topmost section contains general-purpose libraries, while the bottom contains resources specific to clinical text.

Abbreviations: API, application programming interface; BERT, bidirectional encoder representations from transformers; NA, not any; NER, named-entity recognition; NLP, natural language processing; UMLS, Unified Medical Language System.

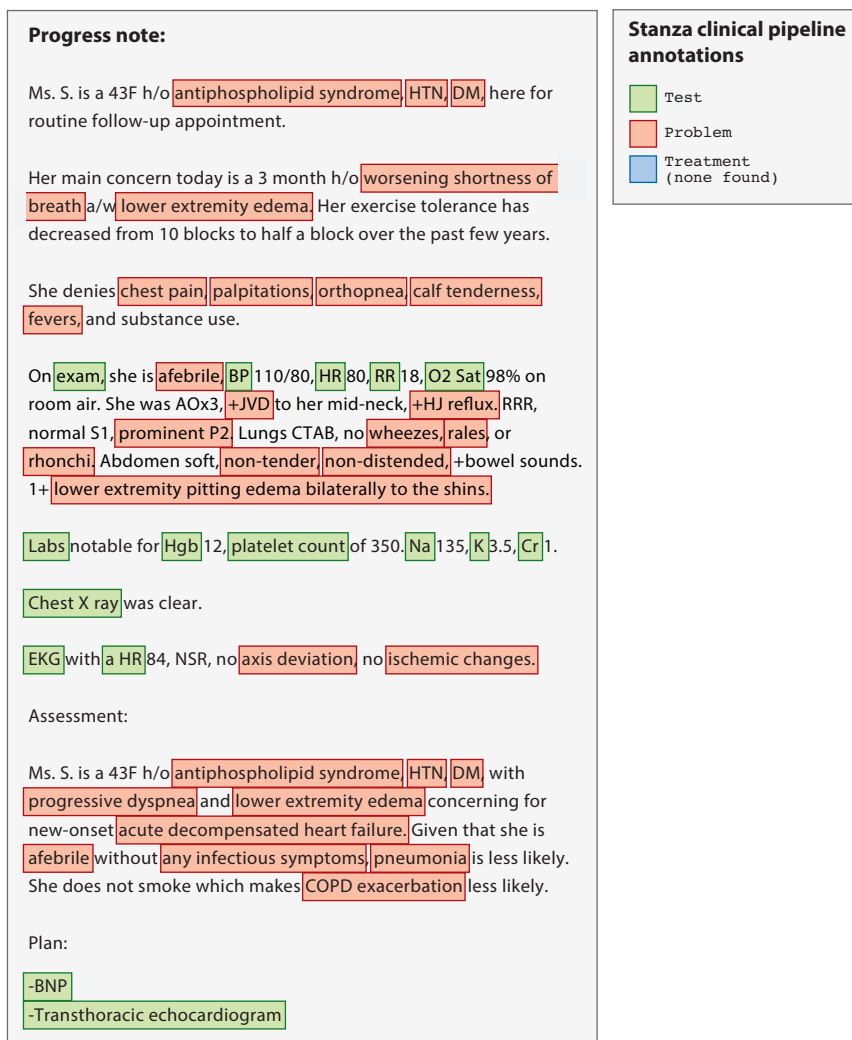


Figure 1

A sample clinical progress note (not a real patient) with named-entity annotations provided by the Stanza clinical text-processing pipeline, trained using data from the 2010 i2b2/VA challenge. The Stanza pipeline tags three types of named entities: **treatment**, **problem**, and **test**. For this particular note, no treatment entities were found. Abbreviations: a/w, along with; AOx3, alert and oriented to person, time, and place; BNP, brain natriuretic peptide test; BP, blood pressure; COPD, chronic obstructive pulmonary disease; Cr, creatinine; CTAB, clear to auscultation bilaterally; DM, diabetes mellitus; Hgb, hemoglobin; HJ reflux, hepatojugular reflux; h/o, history of; HR, heart rate; HTN, hypertension; i2b2/VA, Informatics for Integrating Biology and the Bedside/Veterans Affairs; JVD, jugular vein distention; NSR, normal sinus rhythm; O2 Sat, oxygen saturation; P2, heart sound produced by closure of pulmonic valve; RR, respiratory rate; RRR, regular rate and rhythm (of pulse); S1, heart sound produced by closure of atrioventricular (mitral and tricuspid) valves.

Embedding:

mathematical representation of a word, phrase, or other piece of text, typically a vector of a fixed length, designed so that pieces of text with similar meaning have similar vectors

The simplest NER systems are dictionary based: They simply compare text strings to a list of terms from a specific category, such as disease names. While these approaches are common and frequently yield acceptable performance on clinical text (47, 48), modern clinical NER systems more commonly employ machine learning models adapted for sequence data, including conditional random fields (CRFs), recurrent neural networks (RNNs), and RNN variants such as long short-term memory networks (LSTMs). For example, Stanza (**Figure 1**) uses pretrained character-level language models (49) fed into Bi-LSTM-CRF sequence taggers (32, 50, 51). Trained using corpora hand annotated with entity type(s) of interest, these algorithms learn to identify features of a text string and its surrounding context that predict whether it is one of the desired types.

What constitutes a feature depends on the system. Traditionally, NER algorithms have selected features from predefined sets, including morphological (capitalization and punctuation patterns, presence/absence/location of numbers, etc.), syntactic (parts of speech, grammatical dependencies, etc.), semantic (membership in a lexicon, position in an ontology, etc.), and other specialized or hand-coded features (52). Some systems incorporate pretrained word or phrase embeddings (Section 4), and modern neural network-based NER systems often learn higher-order embeddings directly from patterns in the text itself (50, 53).

A major advantage of machine learning-based NER systems is their flexibility. The same system can often learn to tag different types of entities simply by swapping training datasets. For example, while **Figure 1** shows `test`, `problem`, and `treatment` annotations, the Stanza system has also been trained using a corpus of 150 chest CT (computerized tomography) radiology reports (54) to tag `anatomy`, `anatomy modifier`, `observation`, `observation modifier`, and `uncertainty concepts` (32). The reverse is also true: Different machine learning algorithms can be trained using the same training data. In fact, when looking for an NER system or any other clinical NLP system, a useful strategy is to identify an annotated corpus for that task and look for papers that have cited the corpus. For example, like Stanza, the CliNER (Clinical NER) system (38) was trained using concept annotations from the 2010 i2b2/VA NLP challenge (46). Other widely cited systems trained on the same dataset are the Bi-LSTM-CRF systems by Chalapathy et al. (55) and Unanue et al. (56) and Tang et al.'s system combining support vector machines (SVMs) with CRFs (57).

3.2. Domain Specificity and Key Challenges

It is worth pausing here to consider the conceptual and practical challenges illustrated by **Figure 1**. First, it is clear that NER only makes sense when the entities involved are discrete and have well-defined locations in text. Many clinically important concepts, such as income, housing, and employment history, are unlikely to be described using simple and consistent terminology that can be picked up by NER algorithms. Second, not all of the NER annotations in **Figure 1** are correct or meaningful without consideration of the surrounding context. Labeling the term “acute decompensated heart failure” as a `problem` entity, for example, means little without the qualification that it is a suspected, not definite, diagnosis.

There are also practical concerns regarding NER model training and maintenance. Although dozens of different clinical NER systems have been developed, many are now obsolete, and not all are released as production-ready code (i.e., easy to download and use). In addition, if one is interested in an entity class for which no preannotated corpus or pretrained model is available, there is no alternative but to train one's own system; this means either defining a rules-based approach or creating a custom, annotated training set.

Finally, there is the issue of domain specificity. Clinical text is complex, incorporating specialized medical terms, numerical measures and scores, abbreviations (see **Figure 1** caption),

misspelled words, and poor grammar (20). Many high-quality general domain NER models exist, such as those from the Stanford CoreNLP (25) and spaCy libraries (<https://spacy.io>). However, these are trained using general domain text, such as telephone conversations, newswire, newsgroups, broadcast news, broadcast conversation, and blogs. As such, they tag somewhat generic entities like person, number, and place names, which may or may not be relevant in a clinical text mining context. A second class of systems are those that have been trained using biomedical text, usually from PubMed research articles and abstracts. Recent examples are the scispaCy library (27) and the transformer-based language model BioBERT, fine-tuned for NER (58). These systems often tag entity types that are relevant to clinical text, such as gene names; however, because they were trained using scientific writing, they may suffer reduced accuracy on clinical text. The issue of domain specificity and the need for domain-specific models extend beyond NER, affecting nearly all tasks in clinical text mining.

3.3. Concept Normalization

The output of a clinical NER system (**Figure 1**) is a set of named entities of one or more types. The obvious downside to such output is that it tells one nothing about the entities except their type(s); for example, there is no way of knowing that the strings “HTN” and “hypertension” refer to the same concept—even if they are in the same note and both labeled as `problem`. Likewise, although an NER system may recognize multiword phrases (e.g., “lower extremity pitting edema bilaterally to the shins,” line 13, **Figure 1**), it does not understand how the component words contribute to the meaning of each phrase and it cannot easily connect a given phrase to coreferent phrases, even from the same passage (e.g., “lower extremity edema,” line 5, **Figure 1**).

Concept normalization, also known as entity linking, is the task of assigning a unique identity to each entity name mentioned in a text. In the clinical domain, this typically involves mapping each entity name to a known concept from a structured terminology or ontology. The task is closely related to NER and, indeed, systems often combine the two processes (59). Coreference resolution, in which strings referring to the same entity (e.g., a pronoun and its antecedent) are grouped, is a similar task; it is essentially normalization without the ontology mapping step (60, chapter 22).

Clinical text is incredibly diverse (61), and practitioners from different medical specialties, or who have been trained at different institutions, will often choose different terms for the same concept. The Unified Medical Language System (UMLS), a project begun in 1986 at the National Library of Medicine, was designed to address this issue (36). UMLS is a compendium of biomedical ontologies and terminologies in which concepts occurring across multiple resources are mapped to a single concept unique identifier (CUI). Today, the predominant strategy for clinical concept normalization is to map a given text string to one of these CUIs. End-to-end clinical text mining systems like MedLEE (now Health Fidelity; 10), MetaMap (30), MetaMap Lite (31), and cTAKES (11) all have this functionality. The CLAMP system (62) provides an easy-to-use graphical user interface for building and deploying clinical NLP pipelines, including UMLS mapping.

The same clinical progress note analyzed by Stanza in **Figure 1** is shown in **Figure 2**, this time with annotations produced by cTAKES, a popular system developed at the Mayo Clinic (11). A selection of the 122 detailed UMLS mappings produced by cTAKES is shown in **Table 2**. In addition to UMLS-based concept normalization, cTAKES provides negation detection (63) and can identify uncertainty and experiencer (whether the statement refers to the patient or, e.g., a family member). The results shown in **Figure 2** are from the default cTAKES pipeline, i.e., what one could expect from running cTAKES out of the box. Most of the annotations are correct; for

Coreference

resolution: grouping strings from a passage that refer to the same entity, such as a pronoun and its antecedent

Negation detection:

identifying whether a term or concept is negated in the text; simple pattern-based algorithms, such as NegEx, often suffice

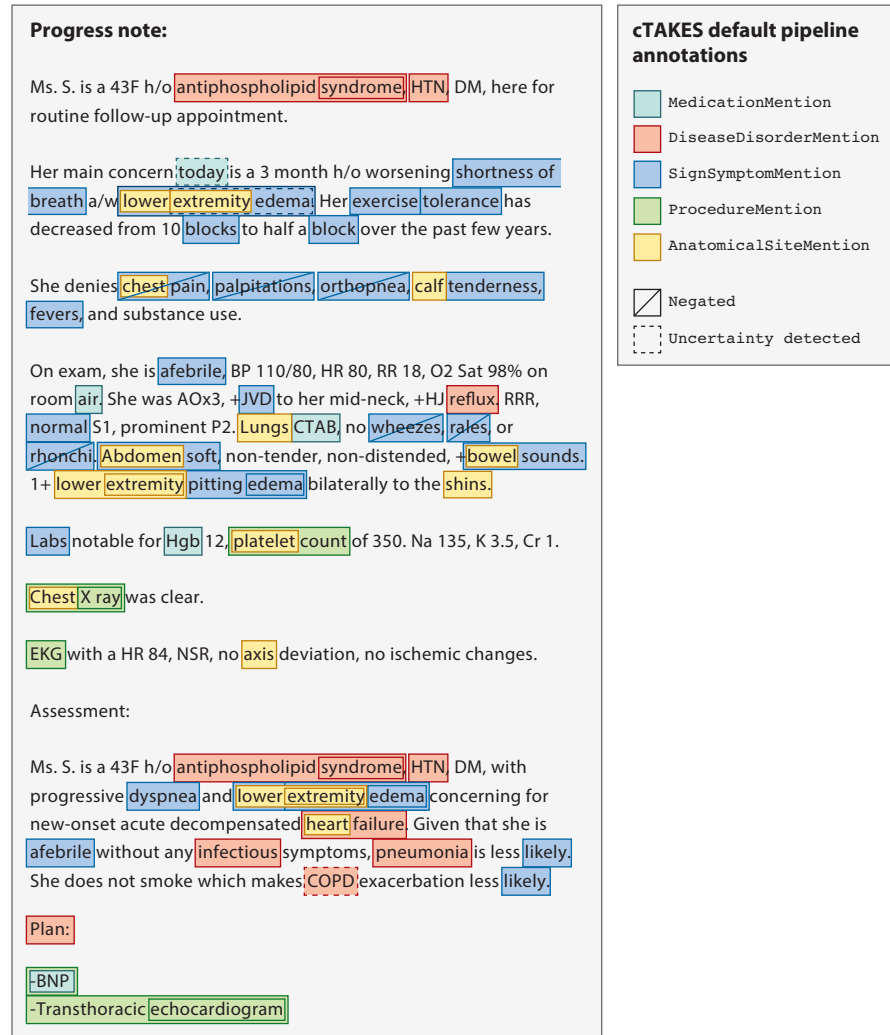


Figure 2

The same clinical progress note as in **Figure 1**, with annotations provided by the cTAKES (version 4.0) default pipeline. The cTAKES pipeline detects negation and uncertainty and maps each entity to its corresponding concept in UMLS. A selection of the UMLS concepts found in this note are listed in **Table 2**. Abbreviations: a/w, along with; AOX3, alert and oriented to person, time, and place; BNP, brain natriuretic peptide test; BP, blood pressure; COPD, chronic obstructive pulmonary disease; Cr, creatinine; CTAB, clear to auscultation bilaterally; DM, diabetes mellitus; Hgb, hemoglobin; HJ reflux, hepatojugular reflux; h/o, history of; HR, heart rate; HTN, hypertension; JVD, jugular vein distention; NSR, normal sinus rhythm; O2 Sat, oxygen saturation; P2, heart sound produced by closure of pulmonic valve; RR, respiratory rate; RRR, regular rate and rhythm (of pulse); S1, heart sound produced by closure of atrioventricular (mitral and tricuspid) valves; UMLS, Unified Medical Language System.

example, cTAKES correctly maps the string “HTN” to the normalized concept *hypertensive disease* (CUI C0020538) and understands that “shortness of breath” is a synonym for *dyspnea* (CUI C0013404). A key shortcoming, however, is cTAKES’ reliance on dictionary-based lookups to identify and normalize named entities. This is apparent in **Figure 2**, where cTAKES labels the

strings “CTAB,” “BNP,” and “Hgb” as medications because of spurious UMLS mappings (e.g., “CTAB” maps to cetrimonium bromide in UMLS). If the specificity of extracted medication terms were crucial for one’s application, it might therefore make sense to include a dedicated NER system for medication names in the cTAKES pipeline. In addition, depending on the application, full concept normalization may not be necessary; in one recent study (64), using cTAKES annotations as features in a 30-day readmission model yielded no better performance than *n*-grams.

Like NER, clinical concept normalization is still an active area of research. For those interested in this task, a good place to start is the disorder normalization systems built for the SHARE/CLEF (Shared Annotated Resources/Conference and Labs of the Evaluation Forum) eHealth 2013 Evaluation Lab, a community NLP challenge focusing on clinical NER and concept normalization (41, 65). DNorm (33, 61) was the top-performing system on the concept normalization task, deploying a pairwise learning-to-rank approach that was the first of its kind in the clinical concept normalization literature. More recent studies have applied deep learning models to the same task and dataset (66, 67).

***n*-gram:** contiguous sequence of *n* items in text; in NLP, this most often refers to a sequence of *n* words, but it can also refer to sequences of characters, syllables, etc.

3.4. Numbers, Ranges, and Sections

There are a few information extraction tasks of particular importance to clinical text for which dedicated systems have been developed. These systems are generally rule based and rely on regular expressions (68). For example, extraction of lab values and vital signs is a distinct task from NER because it requires interpreting numeric values and ranges. The Valx system (69) extracts and structures lab test comparison statements, although so far it has only been applied to trial descriptions from ClinicalTrials.gov. The convolutional neural network (CNN)-based system developed by Xie et al. (70) identifies blood pressure readings, determines the exactness of the readings, and classifies the readings into three classes: general, treatment, and suggestion. Their machine learning-based workflow could be adapted to extract other types of numeric values.

Section identification is another task somewhat unique to the clinical text mining literature. It involves identifying the section labels associated with each span of text within a note (e.g., “Progress Note,” “Assessment,” and “Plan” in **Figure 1**), which informs the interpretation of whatever is found there. To date, the only section identification system used outside the institution in which it was developed is the SecTag system by Denny et al. (71). A complete review of section identification methods and systems can be found in Reference 72.

Table 2 Examples of cTAKES annotations associated with the note in **Figure 2**

Line number(s)	Annotation type	Original string	Normalized term	UMLS concept ID
2	DiseaseDisorderMention	HTN	Hypertensive disease	C0020538
4–5	SignSymptomMention	shortness of breath	Dyspnea	C0013404
7	SignSymptomMention	chest pain (negated)	Chest pain (negated)	C0008031
10	SignSymptomMention	JVD	Jugular venous engorgement	C0425687
16	ProcedureMention	EKG	Electrocardiography	C1623258
24	ProcedureMention	Transthoracic echocardiogram	Transthoracic echocardiography	C0430462
10	DiseaseDisorderMention	reflux	Gastroesophageal reflux disease	C0017168
11	MedicationMention	CTAB	Cetrimonium bromide	C0951233
23	MedicationMention	BNP	Nesiritide	C0054015

The annotations in the topmost section are correct mappings, and those in the bottom section are incorrect mappings. There were 122 unique cTAKES annotations for this note.

Transfer learning:

storing knowledge gained while solving one problem for use on a different, related problem; incorporating pretrained embeddings into task-specific downstream models is a form of transfer learning

4. EMBEDDINGS AND PRETRAINING

The core idea behind concept normalization (Section 3.3) is semantic relatedness; two terms can look different yet describe the same concept. However, semantic relatedness extends beyond the dichotomy of same versus different; terms can have degrees of similarity (e.g., “dog” versus “cat” as opposed to “dog” versus “volcano”) and can be similar in different ways (e.g., “queen” versus “king” as opposed to “queen” versus “president”). Modern NLP systems represent this idea mathematically using a construct called an embedding.

4.1. Word, Phrase, and Character Embeddings

An embedding is a semantically meaningful mathematical representation of a word, phrase, or other piece of text. Usually a vector, it is designed in such a way that words and phrases with similar meanings have similar vectors. Meaning is difficult to represent using numbers, so embedding methods replace “meaning” with “context” and build vectors to reflect usage patterns, typically within large, unlabeled corpora. The NLP subfield of distributional semantics, which originated with Latent Semantic Analysis in 1988 and reached a milestone with the development of word2vec (73) and GloVe (74) in 2013–2014, is a collection of methods all built around the central goal of creating vector space embeddings of words and phrases that reflect how they are used in context. To compare the meaning of two words, one simply calculates the cosine similarity of their corresponding vectors.

From a clinical text mining standpoint, embeddings are useful in two ways. First, because they do not require annotated corpora for training, it is easy to create embeddings that are specific to clinical text or that capture regularities of expression within a particular clinical subfield or institution. These will often outperform general-domain embeddings on clinical text mining tasks (53). Specialized clinical text embeddings have been used to improve clinical NER (75), resolve abbreviations in clinical text (76), expand a structured lexicon of radiology terms (77), and build specialized lexicons from scratch (78). Second, an embedding can incorporate structured information beyond what is found in the text (79), and embeddings have been created to represent CUIs (80), documents (81, 82), or entire patient records (83). Any task in which the notion of similarity is important, particularly when that similarity is based on patterns in text, can probably benefit from embeddings.

For more information about embeddings, readers are encouraged to consult Turney & Pantel (84) for a review of early methods and Kalyan et al. (85) for a review of embedding methods currently in use in clinical text mining.

4.2. Contextual Embeddings and Pretraining

Until the last few years, embeddings consisted of one vector per entity, that is, one vector per word, phrase, or document. However, novel neural network architectures (22) have permitted the creation of embeddings that vary depending on the context; this has expanded the representational power of embedding methods and led to the creation of massive pretrained language models like BERT (bidirectional encoder representations from transformers) (29) and GPT-3 (openai.com). These models are generally too resource intensive to be trained from scratch. Instead, a transfer learning approach (86) is used in which models trained on general-domain corpora are either further pretrained or fine-tuned on clinical text for use in clinical text mining tasks (85). For example, Alsentzer et al. recently trained BERT models on 2 million notes from the MIMIC-III (Medical Information Mart for Intensive Care) (87) database. They produced two models, one for generic clinical text and another for discharge summaries, which they released publicly

(34). They and others have demonstrated that BERT models fine-tuned on clinical corpora improve the state of the art on clinical NER, deidentification, inference, and concept normalization tasks (88, 89), although in at least one case, UMLS features still contributed valuable additional information (90).

The downside of these models is that they require some technical sophistication to adapt and apply. Whereas the original word2vec could be run on a plain text corpus using a single script and output vectors to a text file, using BERT requires knowledge of how to wire up a pretrained model to task-specific output layers for fine-tuning. However, it is likely that end-to-end clinical text processing systems, like cTAKES, will begin to incorporate BERT and related methods into different annotation modules as the technology develops.

5. TEXT CLASSIFICATION

Text classification is perhaps the most sought-after application of clinical text mining. A recent survey (22) found that of 212 clinical text mining papers employing deep learning methods, 88 (41.5%) focused on text classification; text classification and NER together encompassed 75.5% of articles. The goal of text classification is to classify documents (or sentences, phrases, etc.) into two or more discrete categories. Examples from the clinical domain include classifying primary care descriptions of lower back pain into acute versus chronic back pain (91), distinguishing normal versus abnormal knee MRI (magnetic resonance imaging) reports (92), and assessing whether a patient is a current or former smoker versus a nonsmoker based on clinical notes (93). Text classification is a modeling task—typically, it is its own goal. Often it will incorporate features identified through information extraction (Section 3), like named entities or CUIs, or embeddings (Section 4).

A recent systematic review of clinical text classification describes standard text classification algorithms, as well as popular approaches to preprocessing, feature selection, and training set construction (20). An older but still relevant review has surveyed text classification methods for automated clinical coding (94). In general, text classification methods for clinical text are similar to those for other domains, with the exception that specialized medical resources, such as UMLS, often serve as additional sources of features.

5.1. Feature Construction and Selection

The use of individual words or n -grams as features, while common in text classification (95, chapter 13), often results in undesirable levels of feature sparsity when applied to clinical text. As a result, feature selection and dimensionality reduction methods are of particular importance in clinical text classification. Feature selection based on TFIDF (term frequency-inverse document frequency) weighting (95, chapter 6) is common, as are embeddings (Section 4), which turn a potentially unmanageable number of word and text features into dense representations of fixed dimensionality (96). Concept normalization (Section 3.3) also plays a particularly important role in clinical text classification; it is common to preprocess clinical text with a system like cTAKES or MetaMap to merge different term and phrase variants into the same structured concept, and then use those concepts in a classification model (97, 98). It is also possible to exploit parent–child relationships from the UMLS hierarchy to create additional features, e.g., by including all parent terms for a given concept. Such ontology-guided feature engineering has been shown to improve performance on downstream clinical text classification tasks (99). Finally, one can choose a classification algorithm that provides implicit feature selection. In one study, elastic net (100) was used to classify ICU patients into risk strata based on the text of nursing notes. It reduced the number of text features by over 1,000-fold while maintaining near-optimal performance (101).

Weak supervision:

supervised learning using weak (noisy) labels; for example, simple heuristic rules (labeling functions) may be used to create large, weakly annotated training sets

Distant supervision:

supervised learning using training signals that do not directly label the training examples (e.g., using structured clinical data to train a text mining algorithm, since these data are associated with patients or encounters, not sentences or documents)

Relation extraction:

assigning a structured form to a relationship between entities; typically this form includes the normalized entities and a label denoting the nature of their relationship

5.2. Deep Learning for Clinical Text Classification

Aside from those that have employed task-specific rules (Section 2.2), the majority of clinical text classification studies to date have used standard supervised machine learning algorithms, including SVMs, naive Bayes, random forests, and boosting (92, 102, 103). However, over the past five years, deep learning algorithms have begun to displace other classifiers. One of their key advantages is a reduced need for feature engineering; embeddings of words, phrases, and higher-order text structures can be learned as part of the overall training process or incorporated via transfer learning from other pretrained models. Several studies have deployed CNNs with high success on a variety of clinical text classification tasks: assigning diagnosis codes (104, 105), classifying radiology reports (21, 106), subtyping diseases (91), and determining the presence or absence of comorbidities (107). Alternative neural network architectures, such as LSTMs and attention networks, are commonly used in text classification tasks in the general NLP domain, although as of this writing, CNNs have been the dominant architecture in clinical text classification (22, 108). One recent paper exemplifies the end-to-end deep learning approach to clinical text classification, tying rule-based features together with word- and UMLS-based concept embeddings in a single CNN-based classifier (107).

6. WEAK AND DISTANT SUPERVISION

As discussed in Section 2.2, machine learning approaches to clinical NLP generally suffer from a lack of training data (23). In addition, existing clinical information extraction and text classification models have generally been trained using the same few annotated datasets (46, 65, 109–111), which restricts the range and quality of annotations they produce. Most applied clinical text mining projects will therefore confront, at some point, the problem of insufficient or inappropriate training data. Two practical solutions to this problem are weak supervision and distant supervision. Weak supervision is the act of creating silver standard training data by applying a weak, or noisy, labeling function to large amounts of unlabeled data. Distant supervision is a related practice in which external data sources, such as knowledgebases, serve as training signals without labeling the training examples directly. One can view distant supervision as a form of weak supervision, and in practice the terms are often used interchangeably.

The paradigmatic clinical text mining example of distant supervision is using structured information from the EHR, such as ICD (International Classification of Diseases) codes, as a labeling mechanism for unstructured text documents. For example, outcomes such as in-hospital mortality (16), hospital readmission (112, 113), and reportable adverse events (114) are routinely captured in the course of health system operations. Although this information is typically attached to patients or encounters, not individual text documents, one can use it as a source of noisy training labels for discharge summaries or other narrative documents attached to the encounters. These noisy labels then serve as a source of supervision for text classification algorithms. Similar results have been achieved using structured ICD9/10 (ICD, 9th or 10th Revision) diagnosis (115–117) and procedure codes (118) as class labels. However, this technique is somewhat limited to the task of document classification; to obtain labels for specific words or text spans (i.e., for NER or relation extraction), one needs a labeling mechanism that works directly on the text.

An alternative is to apply simple heuristic rules to create noisy labels. For example, Wang et al. used keyword-based weak labels for two separate tasks: smoking status classification and hip fracture classification (93). Importantly, they noted that their best-performing deep learning classifier, a CNN, was robust to the massive label noise created by the weak labeling. Their paper was, to my knowledge, the first to apply a combination of weak supervision and deep learning to clinical text classification; most earlier applications of weak supervision in the biomedical domain

focused on images or text from biomedical research articles. Two earlier studies of note in the biomedical domain are Sabbir et al.'s study of distant supervision for biomedical word sense disambiguation (119) and Fries et al.'s description of the SwellShark system (120), a generative model for biomedical NER that uses lexicons and ontologies for weak labeling. The Snorkel system, on which SwellShark is based, was recently used to weakly label clinical notes for the purposes of extracting implant details and reports of complications and pain after hip replacement; the weakly labeled notes were then used to train deep learning models to recognize pain–anatomy and complication–implant relations (121). These methods improved classification performance by 12.8–53.9% over rule-based methods and detected over six times as many complication events as structured data alone.

Alternative approaches to the efficient annotation of training sets for clinical text mining include crowdsourcing and active learning. Crowdsourcing is not usually a viable option in the clinical domain because of privacy concerns. Active learning is a strategy for minimizing annotation effort by iteratively sampling subsets of data for human annotation based on the current predictions of a supervised learning algorithm (122, 123). However, it still requires recruiting one or more experts to create the annotations.

7. RELATION EXTRACTION AND INFERENCE

Relation extraction is the task of assigning a structured form to a relationship between or among entities based on how they are described in text. Typically this form includes the categories of the involved entities and a label denoting the nature of their relationship, such as “symptom *sign_of* disease” or “test *reveals* problem.” For example, the phrase “progressive dyspnea and lower extremity edema concerning for new-onset acute decompensated heart failure” from the last paragraph in **Figure 1** contains two different “symptom *sign_of* disease” relations. Relation extraction is usually framed as a text classification problem in which sentences or dependency paths (produced by dependency parsing; see margin definition) are classified into groups corresponding to relational labels. It is related to the task of creating a knowledgebase, which represents text as a network of structured relations over which inference can be performed to generate new knowledge (124).

Although ordinarily discussed alongside other information extraction tasks, such as NER, relation extraction is arguably one step closer to true language understanding. NER and text classification simply label text; they do not address compositionality, the combining of individual facts to generate composite ideas. Compositionality presents a particularly important challenge for clinical text mining because clinical writing reflects a high level of assumed knowledge, as well as unstated implications about the temporal and causal ordering of events. Current clinical text mining systems possess no ability to reason, as a human would, about the relationships among laboratory and clinical findings and specific diagnoses or treatments (in **Figure 1**, the meaning of a clear chest X-ray or the implication of pitting edema for a diagnosis of heart failure). Such reasoning requires incorporation of external knowledge derived from, e.g., textbooks or research articles. Relation extraction is a first step in this direction.

7.1. Methods for Clinical Relation Extraction

Modern clinical relation extraction systems are generally based on deep learning models, such as CNNs with pretrained word2vec embeddings (125), segment CNNs (126), and coupled Bi-LSTMs with CNNs incorporating dependency path features (127), or other machine learning methods like SVMs (128, 129). They are typically built and evaluated using annotated corpora, such as the relation extraction corpus from the 2010 i2b2/VA dataset (46), which we have seen earlier; indeed, the five studies just mentioned all used this dataset. The recent 2018 n2c2

Dependency path:

a list of all the edges traversed when moving from one entity to another through a dependency graph; it tends to capture parts of the sentence relevant to the relationship between the two entities

Dependency parsing:

representing the syntactic structure of a sentence as a set of directed, binary grammatical relations between pairs of words (or lemmas); the output is a graph structure called a dependency graph

Compositionality:

a principle from philosophy and mathematical logic, usually attributed to George Boole, stating that the meaning of a complex expression is determined by the meanings of its constituent expressions and the rules used to combine them

Entailment: a type of relation between two segments of text in which one implies the other; that is, the truth of the second statement follows from the first

Natural language inference: also called entailment recognition, the task of determining whether a given hypothesis (statement two) can be inferred from a given premise (statement one)

(National NLP Clinical Challenges) shared task on adverse drug event relations (130) provides a recent snapshot of the field; of the top ten systems, five used deep learning, three used SVMs, one used a random forest, and one used a rule-based algorithm.

One particular relational class that has been the focus of considerable research in recent years is temporal relations, reviewed in detail in Reference 131. A standard language has been developed for annotating temporal relations in text, including events (EVENTs), time expressions (TIMEXs), and relations between EVENTs and TIMEXs (TLINKs). This formalism has led to the creation of two major annotated corpora for clinical temporal relation extraction: the THYME corpus (132) and the 2012 i2b2 temporal relations corpus (110). Methods for temporal relation extraction have followed those developed for other clinical relation extraction tasks; earlier papers used models such as CRFs and SVMs (133), while later papers apply deep learning approaches such as CNNs (134), Bi-LSTMs (135), and BERT (136).

7.2. Inference and Entailment

Natural language inference (NLI) is a close relative of relation extraction with a longstanding presence in NLP, the goal of which is to determine whether one statement (the hypothesis) can be inferred from another (the premise). As of 2018, the clinical NLP community lacked any annotated corpora for NLI, owing in part to the difficulty and expense of getting medical experts to produce annotations and the inability to share patient data with nonexpert (e.g., crowd-worker) annotators. However, Romanov & Shivade (137) recently produced the MedNLI dataset to facilitate NLI research in the clinical domain. Starting with premises from the MIMIC-III (87) dataset, physicians were asked to write sentences that (*a*) were definitely implied by the premise, (*b*) were neither contradicted nor implied by the premise, and (*c*) were definitely contradicted by the premise. Although the task is still in its infancy, shared tasks built around the MedNLI dataset have led to multiple new approaches for NLI in this domain, including BERT-BiLSTM-Attention architectures (138) and state-of-the-art ESIM (enhanced sequential inference model) architectures coupled with knowledge-enhanced word representations based on UMLS (139, 140).

8. CONCLUSION

The volume of EHR data in the United States has skyrocketed in recent years. In 2008, only 42% of office-based physicians reported access to an EHR; this number had climbed to 86% a decade later (145). A favorable policy environment, created by the HITECH (Health Information Technology for Economic and Clinical Health) Act of 2009 and fueled by the 21st Century Cures Act of 2016, has promoted the meaningful use of EHRs and other observational health data to inform patient care, improve health system operations, facilitate research, and provide real-world evidence for FDA (Food and Drug Administration) approval. Over this same time period, methodological advances in machine learning (22, 86), the creation of dedicated clinical text processing software (11, 30), and the seemingly continuous development of high-performing predictive and diagnostic algorithms (6) have fueled enthusiasm about the ability of data and data science to change the way we deliver healthcare.

Amidst such excitement, it would be easy to overlook the fact that most predictive models built on EHR data have focused on outcomes captured in structured data fields, such as mortality, readmissions, length of stay, and diagnosis codes (83, 146). In addition, a recent systematic review found no net performance benefit of more sophisticated machine learning methods over logistic regression in clinical prediction models (147). Both of these observations can be explained, in part, by the limitations of clinical text mining. To date, the vast quantities of text contained within EHRs have primarily been treated as a source of features for downstream learning algorithms,

Table 3 A review of reviews

Year	Author(s)	Title	Reference
2011	Chapman et al.	Overcoming barriers to NLP for clinical text: the role of shared tasks and the need for additional creative solutions	141
2016	Ford et al.	Extracting information from the text of electronic medical records to improve case detection: a systematic review	142
2019	Koleck et al.	Natural language processing of symptoms documented in free-text narratives of electronic health records: a systematic review	5
2017	Kreimeyer et al.	Natural language processing systems for capturing and standardizing unstructured clinical information: a systematic review	8
2019	Khattak et al.	A survey of word embeddings for clinical text	143
2019	Mujtaba et al.	Clinical text classification research trends: systematic literature review and open issues	20
2020	Spasic et al.	Clinical text data in machine learning: systematic review	23
2010	Stanfill et al.	A systematic literature review of automated clinical coding and classification systems	94
2018	Velupillai et al.	Using clinical natural language processing for health outcomes research: overview and actionable suggestions for future advances	144
2018	Wang et al.	Clinical information extraction applications: a literature review	17
2020	Wu et al.	Deep learning in clinical natural language processing: a methodical review	22

The field of clinical text mining has been extensively reviewed in prior articles. The reviews selected here are those I found to be particularly useful surveys of specific research areas or the field in general.

improving predictive performance over structured data alone (142, 148, 149), but not enabling the types of fundamentally new studies that would result if clinical text mining systems could reason about text and incorporate prior knowledge the way a human would. Assessing whether a treatment failed or succeeded for a given patient, for example, is still a nearly impossible task to accomplish using EHR data without manual chart review. Even the most cutting-edge healthcare data science companies still employ human curators to extract this type of information from text. This situation limits both the types of questions we can ask of EHR data and the ability of even the most sophisticated predictive algorithms and causal models to answer them.

Modern clinical text mining systems have accomplished a great deal (see **Table 3**). They can now reliably tag a wide variety of clinically relevant entities in text, map them to standard concepts from lexicons and ontologies, detect negation and uncertainty, and understand the person or people to whom they refer. Given sufficient training data, there are now established system architectures for performing tasks like text classification and relation extraction in the clinical domain. Production-grade clinical text mining systems are in use throughout industry and academia, finding wide application in health outcomes research (144), case detection and phenotyping (142), and automated coding and classification (94). Modern deep learning methods, particularly massive language models like BERT, have recently entered the clinical domain, improving state-of-the-art performance on a variety of clinical NLP tasks and rightfully generating much excitement (86). There remain open questions about the fundamental limitations of these methods to process and understand language (150), and to date the rate of publications describing the application of text

mining to EHR data has not kept pace with the field of EHR data mining as a whole (17). However, the field of clinical text mining is also at an exciting turning point, as it is beginning to pursue questions of inference and logic that cut to the heart of what it means to build intelligent machines.

DISCLOSURE STATEMENT

The author is not aware of any affiliations, memberships, funding, or financial holdings that might be perceived as affecting the objectivity of this review.

ACKNOWLEDGMENTS

Many thanks to Cindy Gao for reference checking and editing, Edwin Yoo for writing the fake outpatient progress note that is included in this review, and Guergana Savova, the creator of cTAKES, for her advice on implementing the cTAKES default pipeline. I would also like to acknowledge Tudor Achim, a contributor to Stack Overflow, who was the source of the clear distinction between weak and distant supervision used in the paper (<https://stackoverflow.com/questions/18944805>). Thanks also to Marshall Pierce and Ben Glicksberg, who reviewed the first draft of the paper and provided helpful feedback on both the writing and the content.

LITERATURE CITED

1. Roberts A. 2017. Language, structure, and reuse in the electronic health record. *AMA J. Ethics* 19:281–88
2. Hersh WR, Weiner MG, Embi PJ, Logan JR, Payne PR, et al. 2013. Caveats for the use of operational electronic health record data in comparative effectiveness research. *Med. Care* 51:S30–37
3. Rosenbloom ST, Denny JC, Xu H, Lorenzi N, Stead WW, Johnson KB. 2011. Data from clinical notes: a perspective on the tension between structure and flexible documentation. *J. Am. Med. Inf. Assoc.* 18:181–86
4. Hatef E, Rouhizadeh M, Tia I, Lasser E, Hill-Briggs F, et al. 2019. Assessing the availability of data on social and behavioral determinants in structured and unstructured electronic health records: a retrospective analysis of a multilevel health care system. *JMIR Med. Inform.* 7:e13802
5. Koleck TA, Dreisbach C, Bourne PE, Bakken S. 2019. Natural language processing of symptoms documented in free-text narratives of electronic health records: a systematic review. *J. Am. Med. Inform. Assoc.* 26:364–79
6. Topol EJ. 2019. High-performance medicine: the convergence of human and artificial intelligence. *Nat. Med.* 25:44–56
7. Rajkomar A, Dean J, Kohane I. 2019. Machine learning in medicine. *N. Engl. J. Med.* 380:1347–58
8. Kreimeyer K, Foster M, Pandey A, Arya N, Halford G, et al. 2017. Natural language processing systems for capturing and standardizing unstructured clinical information: a systematic review. *J. Biomed. Inform.* 73:14–29
9. Dalianis H. 2018. *Clinical Text Mining: Secondary Use of Electronic Patient Records*. Cham, Switz.: Springer Int.
10. Friedman C, Alderson PO, Austin JH, Cimino JJ, Johnson SB. 1994. A general natural-language text processor for clinical radiology. *J. Am. Med. Inform. Assoc.* 1:161–74
11. Savova GK, Masanz JJ, Ogren PV, Zheng J, Sohn S, et al. 2010. Mayo clinical text analysis and knowledge extraction system (cTAKES): architecture, component evaluation and applications. *J. Am. Med. Inform. Assoc.* 17:507–13
12. Guzman B, Metzger I, Aphinyanaphongs Y, Grover H, et al. 2020. Assessment of Amazon Comprehend Medical: medication information extraction. arXiv:2002.00481 [cs.CL]
13. Wei WQ, Teixeira PL, Mo H, Cronin RM, Warner JL, Denny JC. 2016. Combining billing codes, clinical notes, and medications from electronic health records provides superior phenotyping performance. *J. Am. Med. Inform. Assoc.* 23:e20–27

14. Liao KP, Cai T, Savova GK, Murphy SN, Karlson EW, et al. 2015. Development of phenotype algorithms using electronic medical records and incorporating natural language processing. *BMJ* 350:h1885
15. Marafino BJ, Park M, Davies JM, Thombley R, Luft HS, et al. 2018. Validation of prediction models for critical care outcomes using natural language processing of electronic health record data. *JAMA Netw. Open* 1:e185097
16. Weissman GE, Hubbard RA, Ungar LH, Harhay MO, Greene CS, et al. 2018. Inclusion of unstructured clinical text improves early prediction of death or prolonged ICU stay. *Crit. Care Med.* 46:1125–32
17. Wang Y, Wang L, Rastegar-Mojarad M, Moon S, Shen F, et al. 2018. Clinical information extraction applications: a literature review. *J. Biomed. Inform.* 77:34–49
18. Castro SM, Tseytlin E, Medvedeva O, Mitchell K, Visweswaran S, et al. 2017. Automated annotation and classification of BI-RADS assessment from radiology reports. *J. Biomed. Inform.* 69:177–87
19. Patterson O, Hurdle JF. 2011. *Document clustering of clinical narratives: a systematic study of clinical sublanguages.* *AMIA Annu. Symp. Proc.* 2011:1099–107
20. Muftaba G, Shuib L, Idris N, Hoo WL, Raj RG, et al. 2019. Clinical text classification research trends: systematic literature review and open issues. *Expert Syst. Appl.* 116:494–520
21. Shin B, Chokshi FH, Lee T, Choi JD. 2017. *Classification of radiology reports using neural attention models.* In *2017 International Joint Conference on Neural Networks (IJCNN)*, pp. 4363–70. New York: IEEE
22. Wu S, Roberts K, Datta S, Du J, Ji Z, et al. 2020. Deep learning in clinical natural language processing: a methodical review. *J. Am. Med. Inform. Assoc.* 27:457–70
23. Spasic I, Nenadic G. 2020. Clinical text data in machine learning: systematic review. *JMIR Med. Inform.* 8:e17984
24. Bird S, Klein E, Loper E. 2009. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit.* Sebastapol, CA: O'Reilly Media
25. Manning CD, Surdeanu M, Bauer J, Finkel J, Bethard SJ, McClosky D. 2014. The Stanford CoreNLP natural language processing toolkit. In *Association for Computational Linguistics (ACL) System Demonstrations*, ed. K Bontcheva, J Zhu, pp. 55–60. Stroudsburg, PA: Assoc. Comput. Linguist.
26. Qi P, Zhang Y, Zhang Y, Bolton J, Manning CD. 2020. Stanza: a Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, ed. A Celikyilmaz, T-H Wen, pp. 101–8. Stroudsburg, PA: Assoc. Comput. Linguist.
27. Neumann M, King D, Beltagy I, Ammar W. 2019. ScispaCy: fast and robust models for biomedical natural language processing. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, ed. D Demner-Fushman, KB Cohen, S Ananiadou, J Tsujii, pp. 319–27. Stroudsburg, PA: Assoc. Comput. Linguist.
28. Řehůřek R, Sojka P. 2010. Software framework for topic modelling with large corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pp. 45–50. Valletta, Malta: ELRA. <http://is.muni.cz/publication/884893/en>
29. Devlin J, Chang MW, Lee K, Toutanova K. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. arXiv:1810.04805 [cs.CL]
30. Aronson AR, Lang FM. 2010. An overview of MetaMap: historical perspective and recent advances. *J. Am. Med. Inform. Assoc.* 17:229–36
31. Demner-Fushman D, Rogers WJ, Aronson AR. 2017. MetaMap Lite: an evaluation of a new Java implementation of MetaMap. *J. Am. Med. Inform. Assoc.* 24:841–44
32. Zhang Y, Zhang Y, Qi P, Manning CD, Langlotz CP. 2020. Biomedical and clinical English model packages in the Stanza Python NLP library. arXiv:2007.14640 [cs.CL]
33. Leaman R, Islamaj Doğan R, Lu Z. 2013. DNorm: disease name normalization with pairwise learning to rank. *Bioinformatics* 29:2909–17
34. Alsentzer E, Murphy JR, Boag W, Weng WH, Jin D, et al. 2019. Publicly available clinical BERT embeddings. arXiv:1904.03323 [cs.CL]
35. Huang K, Altosaar J, Ranganath R. 2019. ClinicalBERT: Modeling clinical notes and predicting hospital readmission. arXiv:1904.05342 [cs.CL]
36. Bodenreider O. 2004. The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic Acids Res.* 32:D267–70

37. Wu Y, Jiang M, Xu J, Zhi D, Xu H. 2017. Clinical named entity recognition using deep learning models. In *AMIA Annu. Symp. Proc.* 2017:1812–19
38. Boag W, Sergeeva E, Kulshreshtha S, Szolovits P, Rumshisky A, Naumann T. 2018. CliNER 2.0: accessible and accurate clinical concept extraction. arXiv:1803.02245 [cs.CL]
39. Goeuriot L, Suominen H, Kelly L, Miranda-Escalada A, Krallinger M, et al. 2020. Overview of the CLEF eHealth Evaluation Lab 2020. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction: Proceedings of the 11th International Conference of the CLEF Association (CLEF 2020)*, pp. 255–71. Cham, Switz.: Springer
40. Luo YF, Sun W, Rumshisky A. 2019. MCN: a comprehensive corpus for medical concept normalization. *J. Biomed. Inform.* 92:103132
41. Pradhan S, Elhadad N, South BR, Martinez D, Christensen L, et al. 2015. Evaluating the state of the art in disorder recognition and normalization of the clinical narrative. *J. Am. Med. Inform. Assoc.* 22:143–54
42. Yang X, Bian J, Wu Y. 2018. Detecting medications and adverse drug events in clinical notes using recurrent neural networks. In *Proc. Mach. Learn. Res.* 90:1–6
43. Jagannatha AN, Yu H. 2016. Structured prediction models for RNN based sequence labeling in clinical text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 856–65. Stroudsburg, PA: Assoc. Comput. Linguist.
44. Liu Z, Yang M, Wang X, Chen Q, Tang B, et al. 2017. Entity recognition from clinical texts via recurrent neural network. *BMC Med. Inform. Decis. Making* 17:67
45. Dernoncourt F, Lee JY, Uzuner O, Szolovits P. 2017. De-identification of patient notes with recurrent neural networks. *J. Am. Med. Inform. Assoc.* 24:596–606
46. Uzuner Ö, South BR, Shen S, DuVall SL. 2011. 2010 i2b2/va challenge on concepts, assertions, and relations in clinical text. *J. Am. Med. Inform. Assoc.* 18:552–56
47. Jung K, LePendur P, Iyer S, Bauer-Mehren A, Percha B, Shah NH. 2015. Functional evaluation of out-of-the-box text-mining tools for data-mining tasks. *J. Am. Med. Inform. Assoc.* 22:121–31
48. Quimbaya AP, Múnera AS, Rivera RAG, Rodríguez JCD, Velandia OMM, et al. 2016. Named entity recognition over electronic health records through a combined dictionary-based approach. *Procedia Comput. Sci.* 100:55–61
49. Akbik A, Blythe D, Vollgraf R. 2018. Contextual string embeddings for sequence labeling. In *Proceedings of the 27th International Conference on Computational Linguistics*, ed. EM Bender, L Derczynski, P Isabelle, pp. 1638–49. Stroudsburg, PA: Assoc. Comput. Linguist.
50. Lample G, Ballesteros M, Subramanian S, Kawakami K, Dyer C. 2016. Neural architectures for named entity recognition. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, ed. K Knight, A Nenkova, O Rambow, pp. 260–70. Stroudsburg, PA: Assoc. Comput. Linguist.
51. Huang Z, Xu W, Yu K. 2015. Bidirectional LSTM-CRF models for sequence tagging. arXiv:1508.01991 [cs.CL]
52. Settles B. 2004. Biomedical named entity recognition using conditional random fields and rich feature sets. In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications (NLPBA/BioNLP)*, pp. 107–10. Stroudsburg, PA: Assoc. Comput. Linguist.
53. Wang Y, Liu S, Afzal N, Rastegar-Mojarad M, Wang L, et al. 2018. A comparison of word embeddings for the biomedical natural language processing. *J. Biomed. Inform.* 87:12–20
54. Hassanpour S, Langlotz CP. 2016. Information extraction from multi-institutional radiology reports. *Artif. Intel. Med.* 66:29–39
55. Chalapathy R, Borzeshi EZ, Piccardi M. 2016. Bidirectional LSTM-CRF for clinical concept extraction. arXiv:1611.08373 [cs.CL]
56. Unanue JJ, Borzeshi EZ, Piccardi M. 2017. Recurrent neural networks with specialized word embeddings for health-domain named-entity recognition. *J. Biomed. Inform.* 76:102–9
57. Tang B, Cao H, Wu Y, Jiang M, Xu H. 2012. Clinical entity recognition using structural support vector machines with rich features. In *Proceedings of the ACM Sixth International Workshop on Data and Text Mining in Biomedical Informatics*, pp. 13–20. New York: Assoc. Comput. Linguist.
58. Lee J, Yoon W, Kim S, Kim D, Kim S, et al. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36:1234–40

59. Luo G, Huang X, Lin CY, Nie Z. 2015. Joint entity recognition and disambiguation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, ed. L Màrquez, C Callison-Burch, J Su, pp. 879–88. Stroudsburg, PA: Assoc. Comput. Linguist.
60. Jurafsky D, Martin JH. 2019. *Speech and Language Processing*. Book Draft, 3rd ed. <https://web.stanford.edu/~jurafsky/slp3>
61. Leaman R, Khare R, Lu Z. 2015. Challenges in clinical natural language processing for automated disorder normalization. *J. Biomed. Inform.* 57:28–37
62. Soysal E, Wang J, Jiang M, Wu Y, Pakhomov S, et al. 2018. CLAMP: a toolkit for efficiently building customized clinical natural language processing pipelines. *J. Am. Med. Inform. Assoc.* 25:331–36
63. Chapman WW, Bridewell W, Hanbury P, Cooper GF, Buchanan BG. 2001. A simple algorithm for identifying negated findings and diseases in discharge summaries. *J. Biomed. Inform.* 34:301–10
64. Afshar M, Dligach D, Sharma B, Cai X, Boyda J, et al. 2019. Development and application of a high throughput natural language processing architecture to convert all clinical documents in a clinical data warehouse into standardized medical vocabularies. *J. Am. Med. Inform. Assoc.* 26:1364–69
65. Suominen H, Salanterä S, Velupillai S, Chapman WW, Savova G, et al. 2013. Overview of the ShARc/CLEF eHealth Evaluation Lab 2013. In *International Conference of the Cross-Language Evaluation Forum for European Languages*, ed. P Forner, H Muller, R Paredes, P Rosso, B Stein, pp. 212–31. Berlin: Springer
66. Li H, Chen Q, Tang B, Wang X, Xu H, et al. 2017. CNN-based ranking for biomedical entity normalization. *BMC Bioinform.* 18:79–86
67. Luo YF, Sun W, Rumshisky A. 2019. A hybrid normalization method for medical concepts in clinical narrative using semantic matching. *AMIA Summits Transl. Sci. Proc.* 2019:732–40
68. Turchin A, Kolatkar NS, Grant RW, Makhni EC, Pendergrass ML, Einbinder JS. 2006. Using regular expressions to abstract blood pressure and treatment intensification information from the text of physician notes. *J. Am. Med. Inform. Assoc.* 13:691–95
69. Hao T, Liu H, Weng C. 2016. Valx: a system for extracting and structuring numeric lab test comparison statements from text. *Methods Inform. Med.* 55:266–75
70. Xie T, Zhen Y, Tavakoli M, Hundley G, Ge Y. 2020. A deep-learning based system for accurate extraction of blood pressure data in clinical narratives. *AMIA Summits Transl. Sci. Proc.* 2020:703–9
71. Denny JC, Spickard A 3rd, Johnson KB, Peterson NB, Peterson JF, Miller RA. 2009. Evaluation of a method to identify and categorize section headers in clinical documents. *J. Am. Med. Inform. Assoc.* 16:806–15
72. Pomares-Quimbaya A, Kreuzthaler M, Schulz S. 2019. Current approaches to identify sections within clinical narratives from electronic health records: a systematic review. *BMC Med. Res. Methodol.* 19:155
73. Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, ed. CJC Burges, L Bottou, M Welling, Z Ghahramani, KQ Weinberger, pp. 3111–19. <https://papers.nips.cc/paper/2013/hash/9aa42b31882ec039965f3c4923ce901b-Abstract.html>
74. Pennington J, Socher R, Manning CD. 2014. GloVe: global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, ed. A Moschitti, B Pang, W Daelemans, pp. 1532–43. Stroudsburg, PA: Assoc. Comput. Linguist.
75. Wu Y, Xu J, Jiang M, Zhang Y, Xu H. 2015. A study of neural word embeddings for named entity recognition in clinical text. In *AMIA Annu. Symp. Proc.* 2015:1326–33
76. Wu Y, Xu J, Zhang Y, Xu H. 2015. Clinical abbreviation disambiguation using neural word embeddings. In *Proceedings of BioNLP 15*, ed. KB Cohen, D Demner-Fushman, S Ananiadou, J Tsujii, pp. 171–76. Stroudsburg, PA: Assoc. Comput. Linguist.
77. Percha B, Zhang Y, Bozkurt S, Rubin D, Altman RB, Langlotz CP. 2018. Expanding a radiology lexicon using contextual patterns in radiology reports. *J. Am. Med. Inform. Assoc.* 25:679–85
78. Fan Y, Pakhomov S, McEwan R, Zhao W, Lindemann E, Zhang R. 2019. Using word embeddings to expand terminology of dietary supplements on clinical notes. *JAMIA Open* 2:246–53
79. Lastra-Díaz JJ, Goikoetxea J, Taieb MAH, García-Serrano A, Aouicha MB, Agirre E. 2019. A reproducible survey on word embeddings and ontology-based methods for word similarity: linear combinations outperform the state of the art. *Eng. Appl. Artif. Intel.* 85:645–65

80. Beam A, Kompa B, Schmaltz A, Fried I, Weber G, et al. 2020. Clinical concept embeddings learned from massive sources of multimodal medical data. *Pac. Symp. Biocomput.* 25:295–306
81. Baume T, Nassour-Kassis J, Cohen R, Elhadad M, Elhadad N. 2017. Multi-label classification of patient notes a case study on ICD code assignment. arXiv:1709.09587 [cs.CL]
82. Banerjee I, Chen MC, Lungren MP, Rubin DL. 2018. Radiology report annotation using intelligent word embeddings: applied to multi-institutional chest CT cohort. *J. Biomed. Inform.* 77:11–20
83. Miotto R, Li L, Kidd BA, Dudley JT. 2016. Deep patient: an unsupervised representation to predict the future of patients from the electronic health records. *Sci. Rep.* 6:26094
84. Turney PD, Pantel P. 2010. From frequency to meaning: vector space models of semantics. *J. Artif. Intel. Res.* 37:141–88
85. Kalyan KS, Sangeetha S. 2020. SECNLP: a survey of embeddings in clinical natural language processing. *J. Biomed. Inform.* 101:103323
86. Peng Y, Yan S, Lu Z. 2019. Transfer learning in biomedical natural language processing: an evaluation of BERT and ELMo on ten benchmarking datasets. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, ed. D Demner-Fushman, KB Cohen, S Ananiadou, J Tsujii, pp. 58–65. Stroudsburg, PA: Assoc. Comput. Linguist.
87. Johnson AE, Pollard TJ, Shen L, Li-Wei HL, Feng M, et al. 2016. MIMIC-III, a freely accessible critical care database. *Sci. Data* 3:160035
88. Si Y, Wang J, Xu H, Roberts K. 2019. Enhancing clinical concept extraction with contextual embeddings. *J. Am. Med. Inform. Assoc.* 26:1297–304
89. Li F, Jin Y, Liu W, Rawat BPS, Cai P, Yu H. 2019. Fine-tuning bidirectional encoder representations from transformers (BERT)-based models on large-scale electronic health record notes: an empirical study. *JMIR Med. Inform.* 7:e14830
90. Xu D, Gopale M, Zhang J, Brown K, Begoli E, Bethard S. 2020. Unified medical language system resources improve sieve-based generation and bidirectional encoder representations from transformers (BERT)-based ranking for concept normalization. *J. Am. Med. Inform. Assoc.* 27(10):1510–19
91. Miotto R, Percha BL, Glicksberg BS, Lee HC, Cruz L, et al. 2020. Identifying acute low back pain episodes in primary care practice from clinical notes: observational study. *JMIR Med. Inform.* 8:e16878
92. Hassanpour S, Langlotz CP, Amrhein TJ, Befera NT, Lungren MP. 2017. Performance of a machine learning classifier of knee MRI reports in two large academic radiology practices: a tool to estimate diagnostic yield. *Am. J. Roentgenol.* 208:750–53
93. Wang Y, Sohn S, Liu S, Shen F, Wang L, et al. 2019. A clinical text classification paradigm using weak supervision and deep representation. *BMC Med. Informat. Decis. Making* 19:1
94. Stanfill MH, Williams M, Fenton SH, Jenders RA, Hersh WR. 2010. A systematic literature review of automated clinical coding and classification systems. *J. Am. Med. Inform. Assoc.* 17:646–51
95. Manning CD, Schütze H, Raghavan P. 2008. *Introduction to Information Retrieval*. Cambridge, UK: Cambridge Univ. Press
96. Shao Y, Taylor S, Marshall N, Morioka C, Zeng-Treitler Q. 2018. Clinical text classification with word embedding features versus bag-of-words features. In *2018 IEEE International Conference on Big Data (Big Data)*, pp. 2874–78. New York: IEEE
97. Buchan K, Filannino M, Uzuner Ö. 2017. Automatic prediction of coronary artery disease from clinical narratives. *J. Biomed. Inform.* 72:23–32
98. Kocbek S, Cavedon L, Martinez D, Bain C, Mac Manus C, et al. 2016. Text mining electronic hospital records to automatically classify admissions against disease: measuring the impact of linking data sources. *J. Biomed. Inform.* 64:158–67
99. Garla VN, Brandt C. 2012. Ontology-guided feature engineering for clinical text classification. *J. Biomed. Inform.* 45:992–98
100. Zou H, Hastie T. 2005. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. B* 67:301–20
101. Marafino BJ, Boscardin WJ, Dudley RA. 2015. Efficient and sparse feature selection for biomedical text classification via the elastic net: application to ICU risk stratification from nursing notes. *J. Biomed. Inform.* 54:114–20

102. Lucini FR, Fogliatto FS, da Silveira GJ, Neyeloff JL, Anzanello MJ, et al. 2017. Text mining approach to predict hospital admissions using early medical records from the emergency department. *Int. J. Med. Inform.* 100:1–8
103. Kavuluru R, Rios A, Lu Y. 2015. An empirical evaluation of supervised learning approaches in assigning diagnosis codes to electronic medical records. *Artif. Intel. Med.* 65:155–66
104. Rios A, Kavuluru R. 2015. Convolutional neural networks for biomedical text classification: application in indexing biomedical articles. In *Proceedings of the 6th ACM Conference on Bioinformatics, Computational Biology and Health Informatics*, pp. 258–67. New York: Assoc. Comput. Mach.
105. Mullenbach J, Wiegrefe S, Duke J, Sun J, Eisenstein J. 2018. Explainable prediction of medical codes from clinical text. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Vol. 1, ed. M Walker, H Ji, A Stent, pp. 1101–11. Stroudsburg, PA: Assoc. Comput. Linguist.
106. Chen MC, Ball RL, Yang L, Moradzadeh N, Chapman BE, et al. 2018. Deep learning to classify radiology free-text reports. *Radiology* 286:845–52
107. Yao L, Mao C, Luo Y. 2019. Clinical text classification with rule-based features and knowledge-guided convolutional neural networks. *BMC Med. Inform. Decis. Making* 19:71
108. Gehrmann S, Dernoncourt F, Li Y, Carlson ET, Wu JT, et al. 2018. Comparing deep learning and concept extraction based methods for patient phenotyping from clinical narratives. *PLOS ONE* 13:e0192360
109. Uzuner O, Bodnari A, Shen S, Forbush T, Pestian J, South BR. 2012. Evaluating the state of the art in coreference resolution for electronic medical records. *J. Am. Med. Inform. Assoc.* 19:786–91
110. Sun W, Rumshisky A, Uzuner O. 2013. Evaluating temporal relations in clinical text: 2012 i2b2 challenge. *J. Am. Med. Inform. Assoc.* 20:806–13
111. Bethard S, Derczynski L, Savova G, Pustejovsky J, Verhagen M. 2015. Semeval-2015 task 6: clinical tempeval. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, ed. P Nakov, T Zesch, D Cer, D Jurgens, pp. 806–14. Stroudsburg, PA: Assoc. Comput. Linguist.
112. Rumshisky A, Ghassemi M, Naumann T, Szolovits P, Castro V, et al. 2016. Predicting early psychiatric readmission with natural language processing of narrative discharge summaries. *Transl. Psychiatry* 6:e921
113. Agarwal A, Baechle C, Behara R, Zhu X. 2017. A natural language processing framework for assessing hospital readmissions for patients with COPD. *IEEE J. Biomed. Health Inform.* 22:588–96
114. Young IJB, Luz S, Lone N. 2019. A systematic review of natural language processing for classification tasks in the field of incident reporting and adverse event analysis. *Int. J. Med. Inform.* 132:103971
115. Osborne JD, Wyatt M, Westfall AO, Willig J, Bethard S, Gordon G. 2016. Efficient identification of nationally mandated reportable cancer cases using natural language processing and machine learning. *J. Am. Med. Inform. Assoc.* 23:1077–84
116. Chen W, Huang Y, Boyle B, Lin S. 2016. The utility of including pathology reports in improving the computational identification of patients. *J. Pathol. Inform.* 7:46
117. Venkataraman GR, Pineda AL, Bear Don't Walk IV OJ, Zehnder AM, Ayyar S, et al. 2020. FasTag: automatic text classification of unstructured medical narratives. *PLOS ONE* 15:e0234647
118. Roysden N, Wright A. 2015. Predicting health care utilization after behavioral health referral using natural language processing and machine learning. *AMIA Annu. Symp. Proc.* 2015:2063–72
119. Sabbir A, Jimeno-Yepes A, Kavuluru R. 2017. Knowledge-based biomedical word sense disambiguation with neural concept embeddings. In *2017 IEEE 17th International Conference on Bioinformatics and Biengineering (BIBE)*, pp. 163–70. New York: IEEE
120. Fries J, Wu S, Ratner A, Ré C. 2017. SwellShark: a generative model for biomedical named entity recognition without labeled data. arXiv:1704.06360 [cs.CL]
121. Callahan A, Fries JA, Ré C, Huddleston JI, Giori NJ, et al. 2019. Medical device surveillance with electronic health records. *NPJ Digital Med.* 2:94
122. Chen Y, Lasko TA, Mei Q, Denny JC, Xu H. 2015. A study of active learning methods for named entity recognition in clinical text. *J. Biomed. Inform.* 58:11–18
123. Kholghi M, Sitbon L, Zuccon G, Nguyen A. 2016. Active learning: a step towards automating medical concept extraction. *J. Am. Med. Inform. Assoc.* 23:289–96

124. Meystre SM, Thibault J, Shen S, Hurdle JF, South BR. 2010. Texttractor: a hybrid system for medications and reason for their prescription extraction from clinical text documents. *J. Am. Med. Inform. Assoc.* 17:559–62
125. Sahu S, Anand A, Oruganty K, Gattu M. 2016. Relation extraction from clinical texts using domain invariant convolutional neural network. In *Proceedings of the 15th Workshop on Biomedical Natural Language Processing*, ed. KB Cohen, D Demner-Fushman, S Ananiadou, J Tsujii, pp. 206–15. Stroudsburg, PA: Assoc. Comput. Linguist.
126. Luo Y, Cheng Y, Uzuner Ö, Szolovits P, Starren J. 2018. Segment convolutional neural networks (Seg-CNNs) for classifying relations in clinical notes. *J. Am. Med. Inform. Assoc.* 25:93–98
127. Li Z, Yang Z, Shen C, Xu J, Zhang Y, Xu H. 2019. Integrating shortest dependency path and sentence sequence into a deep learning framework for relation extraction in clinical text. *BMC Med. Inform. Decis. Making* 19:22
128. Rink B, Harabagiu S, Roberts K. 2011. Automatic extraction of relations between medical concepts in clinical texts. *J. Am. Med. Inform. Assoc.* 18:594–600
129. Munkhdalai T, Liu F, Yu H. 2018. Clinical relation extraction toward drug safety surveillance using electronic health record narratives: classical learning versus deep learning. *JMIR Public Health Surveill.* 4:e29
130. Henry S, Buchan K, Filannino M, Stubbs A, Uzuner O. 2020. 2018 n2c2 shared task on adverse drug events and medication extraction in electronic health records. *J. Am. Med. Inform. Assoc.* 27:3–12
131. Alfattni G, Peek N, Nenadic G. 2020. Extraction of temporal relations from clinical free text: a systematic review of current approaches. *J. Biomed. Inform.* 108:103488
132. Styler WF IV, Bethard S, Finan S, Palmer M, Pradhan S, et al. 2014. Temporal annotation in the clinical domain. *Trans. Assoc. Comput. Linguist.* 2:143–54
133. Tang B, Wu Y, Jiang M, Chen Y, Denny JC, Xu H. 2013. A hybrid system for temporal information extraction from clinical text. *J. Am. Med. Inform. Assoc.* 20:828–35
134. Dligach D, Miller T, Lin C, Bethard S, Savova G. 2017. Neural temporal relation extraction. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, Vol. 2, ed. M Lapata, P Blunsom, A Koller, pp. 746–51. Stroudsburg, PA: Assoc. Comput. Linguist.
135. Tourille J, Ferret O, Neveol A, Tannier X. 2017. Neural architecture for temporal relation extraction: a Bi-LSTM approach for detecting narrative containers. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, ed. R Barzilay, M-Y Kan, Vol. 2, pp. 224–30. Stroudsburg, PA: Assoc. Comput. Linguist.
136. Lin C, Miller T, Dligach D, Bethard S, Savova G. 2019. A BERT-based universal model for both within- and cross-sentence clinical temporal relation extraction. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, ed. A Rumshisky, K Roberts, S Bethard, T Naumann, pp. 65–71. Stroudsburg, PA: Assoc. Comput. Linguist.
137. Romanov A, Shivade C. 2018. Lessons from natural language inference in the clinical domain. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, ed. E Riloff, D Chiang, J Hockenmaier, J Tsujii, pp. 1586–96. Stroudsburg, PA: Assoc. Comput. Linguist.
138. Lee LH, Lu Y, Chen PH, Lee PL, Shyu KK. 2019. NCUEE at MEDIQA 2019: medical text inference using ensemble BERT-BiLSTM-attention model. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, ed. D Demner-Fushman, KB Cohen, S Ananiadou, J Tsujii, pp. 528–32. Stroudsburg, PA: Assoc. Comput. Linguist.
139. Lu M, Fang Y, Yan F, Li M. 2019. Incorporating domain knowledge into natural language inference on clinical texts. *IEEE Access* 7:57623–32
140. Sharma S, Santra B, Jana A, Santosh T, Ganguly N, Goyal P. 2019. Incorporating domain knowledge into medical NLI using knowledge graphs. arXiv:1909.00160 [cs.CL]
141. Chapman WW, Nadkarni PM, Hirschman L, D’avolio LW, Savova GK, Uzuner O. 2011. Overcoming barriers to NLP for clinical text: the role of shared tasks and the need for additional creative solutions. *J. Am. Med. Inform. Assoc.* 18(5):540–43
142. Ford E, Carroll JA, Smith HE, Scott D, Cassell JA. 2016. Extracting information from the text of electronic medical records to improve case detection: a systematic review. *J. Am. Med. Inform. Assoc.* 23:1007–15

143. Khattak FK, Jeblee S, Pou-Prom C, Abdalla M, Meaney C, Rudzicz F. 2019. A survey of word embeddings for clinical text. *J. Biomed. Inform.* X 4:100057
144. Velupillai S, Suominen H, Liakata M, Roberts A, Shah AD, et al. 2018. Using clinical natural language processing for health outcomes research: overview and actionable suggestions for future advances. *J. Biomed. Inform.* 88:11–19
145. Off. Natl. Coord. Health Inf. Technol. 2019. Office-based physician electronic health record adoption. Tech. Rep. Quick-Stat #50, Off. Natl. Coord. Health Inf. Technol., U.S. Dep. Health Hum. Serv., Washington, DC
146. Rajkomar A, Oren E, Chen K, Dai AM, Hajaj N, et al. 2018. Scalable and accurate deep learning with electronic health records. *NPJ Digit. Med.* 1:18
147. Christodoulou E, Ma J, Collins GS, Steyerberg EW, Verbakel JY, Van Calster B. 2019. A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *J. Clin. Epidemiol.* 110:12–22
148. Castro VM, Minnier J, Murphy SN, Kohane I, Churchill SE, et al. 2015. Validation of electronic health record phenotyping of bipolar disorder cases and controls. *Am. J. Psychiatry* 172:363–72
149. Hoogendoorn M, Szolovits P, Moons LM, Numans ME. 2016. Utilizing uncoded consultation notes from electronic medical records for predictive modeling of colorectal cancer. *Artif. Intell. Med.* 69:53–61
150. Marcus G. 2018. Deep learning: a critical appraisal. arXiv:1801.00631 [cs.AI]