

Clinical Natural Language Processing

Modern Clinical Data Science

May 13, 2021



Mount
Sinai

Outline of the Class

1. Introduction
2. Defining the Task and Approach
3. Clinical Information Extraction
4. Embeddings and Pretraining
5. Text Classification
6. Additional Topics: Weak Supervision and NLI

Part I

Introduction

What is NLP?

(From Wikipedia) A subfield of linguistics, computer science, information engineering, and artificial intelligence concerned with... how to program computers to process and analyze large amounts of natural language data.

- Automatic speech recognition
- CCG
- Common sense
- Constituency parsing
- Coreference resolution
- Dependency parsing
- Dialogue
- Domain adaptation
- Entity linking
- Grammatical error correction
- Information extraction
- Language modeling
- Lexical normalization
- Machine translation
- Missing elements
- Multi-task learning
- Multi-modal
- Named entity recognition
- Natural language inference
- Part-of-speech tagging
- Question answering
- Relation prediction
- Relationship extraction
- Semantic textual similarity
- Semantic parsing
- Semantic role labeling
- Sentiment analysis
- Shallow syntax
- Simplification
- Intent Detection and Slot Filling
- Stance detection
- Summarization
- Taxonomy learning
- Temporal processing
- Text classification
- Word sense disambiguation

<https://github.com/sebastianruder/NLP-progress>

Why “NLP” clinical text?

Most of the information we want to use as features or outcomes in clinical predictive models is found in the text of clinical notes (e.g., laboratory, pathology, or progress notes).

- Treatment goals and outcomes (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)
- Symptoms, symptom changes, and their interpretation
- Past medical history and family history
- Primary and secondary complaints
- Psychiatric evaluations and records of therapy sessions
- 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
- Discharge summaries and follow-up plans

Why “NLP” clinical text? (Example from i2b2)

note-88.txt *

1 | 995247940 | ABOHS | 01328259 | | 983241 | 8/25/1998 12:00:00 AM | CEREBRO VASCULAR
ACCIDENT | Signed | DIS | Admission Date: 4/10/1998 Report Status: Signed

2
3 Discharge Date: 1/19/1998
4 PRINCIPAL DIAGNOSIS: LEFT SUBCORTICAL STROKE.
5 PROBLEM LIST: 1) 25 REVIEW OF SYSTEMS: She had no fevers , chills , vomiting or loss
6 3) HYPERTENSION. 26 of consciousness. No seizure activity. No
7 5) STATUS POST T 27 head or neck trauma. No vertigo or tenderness. No recent change
8 CERVICAL DYSPLASI 28 in medications. No chest pain or shortness of breath , no
9 SITES OF ARTHRITI 29 photophobia or neck stiffness.
10 8) VITILIGO. 9)
11 HISTORY OF PRESEN 30 PAST MEDICAL HISTORY: History of diabetes mellitus .
12 presentation , wh 31 goiter , hypertension .
13 return to bed , t 32 trait , status post .
14 and right arm clu 33 cervical dysplasia .
15 all day long. She 34 complaints , just .
16 She has remained 35 hemorrhoids.
17 prior to admisso 36 MEDICATIONS: Hydrochlorothiazide .
18 admission. This e 37 Norvasc 5 milligrams .
19 paralysis (she h 38 once a day , Premarin .
20 waited for it to 39 insulin 12 units .
21 greater than 24 h 40 14 units of regular .
22 had not taken ins 41 also been taking .
23 intake secondary 42 pain. She is also .
24 difficulty swallow 43 once a day.
44 SOCIAL HISTORY: The patient is a widow. No clear history of alcohol use. She has two children .
45 a widow. No clear 46 use. She has two .
47 mobility.

48 FAMILY HISTORY: Familial intermittent paralytic syndrome.
49 PHYSICAL EXAMINATION: Temperature was 97.1 , heart rate 110 ,
50 respirations 20 , blood pressure 130/70 ,
51 oxygen saturation 100% on room air. Glucose is 388. Mental status
52 examination: The patient was alert and oriented times three
53 following three step command across the midline. No apraxia.
54 Positive dysarthria. Repetition intact. Memory 3/3 immediately
55 and 3/3 at 5 minutes. Cranial nerves: Pupils 3-2 bilaterally.
56 Right central facial paralysis. No pallor or droop. Decreased
57 sensation to light touch , pinprick and temperature on the right.
58 Tears well bilaterally. She cannot shrug her right shoulder.
59 Motor and tone flaccid on the right. Strength: Right hemiplegia
60 (can only wiggle toes) , left strength 4-5/5 (poor effort). Deep
61 tendon reflexes are 2+ at the biceps , 1+ at the knees , 0-1 at the
62 ankle with an upgoing right toe. Sensation: Decreased sensation
63 to light touch , pinprick and vibration. Temperature on the right
64 including trunk. Gait not assessed. HEENT: Atraumatic. Neck
65 supple. Lungs clear to auscultation bilaterally. Cardiovascular:
66 Regular rate and rhythm. Abdomen: Soft and nontender. Positive
67 bowel sounds. Extremities: No edema.
68 EKG: Sinus tachycardia at 104 per report when compared with EKG
69 of January , no significant change. Chest x-ray: No
70 evidence of active cardiopulmonary disease. CT scan of the head on
71 8/17/98 showed low density area present in the posterior limb of
72 the left intracapsular capsule consistent with an infarct.
73 Otherwise , no significant abnormality was seen. Mild age
74 appropriate involitional changes were present.
75 LABORATORY: Glucose 388 mg/dL BUN 25 mg/dL Cr 1.2 mg/dL

Unique Features of Biomedical/Clinical Text

- Specialized vocabulary / abbreviations
- Sentence fragments (esp. clinical)
- Long, complex sentences with multiple clauses (esp. scientific)
- High degree of assumed knowledge

It's important to remember that

NLP is not “solved”.

Consider Google search, Siri, Alexa, etc.

Biomedical and Clinical Corpora

- MIMIC III
- i2b2 (now n2c2)
- PubMed
- PubMed Central
- BioCreative
- GENIA corpus
- Others (see Wikipedia “Biomedical Text Mining” article)

The screenshot shows the DBMI Data Portal homepage. At the top, there is a navigation bar with links for "DBMI Data Portal", "Home", "Data Sets", "Data Challenges", "Software", "Contact", and an email link "bipercha@gmail.com". Below the navigation bar, the main content area features a section titled "n2c2 NLP Research Data Sets" with a sub-section "Unstructured notes from the Research Patient Data Repository at Partners Healthcare". A "Need help? Contact us!" button is present. On the right side, the Harvard Biostatistics Institute logo is displayed. The central part of the page contains a "Description" box with text about the Clinical Natural Language Processing (NLG) data sets and a list of challenges from 2006 to 2018. To the right of this box is a vertical column of download links for various challenges, each with a "..." button.

The screenshot shows the PhysioNet website. At the top, there is a navigation bar with links for "Database", "Open Access", "Find", "Share", "About", "News", "Projects", and a user account link "bipercha". A search bar is also present. The main content area features a section titled "MIMIC-III Clinical Database Demo" by "Alistair Johnson", "Tom Pollard", and "Roger Mark". It includes a publication date of "Published: April 24, 2019. Version: 1.4". The page is divided into several sections: "When using this resource, please cite:", "Additionally, please cite the original publication:", "Please include the standard citation for PhysioNet:", "Abstract", "Contents", "Share" (with social media icons), "Access", and "Access Policy".

Mount Sinai Data Warehouse (MSDW) (slide courtesy of Patricia Kovatch)

9.3 million Number of patient records in MSDW	4.5 million Number of patients with clinical data in MSDW collected since 2003	80 million Number of patient visits recorded, including: inpatient, outpatient and ED	107 million Number of ICD-9 and ICD-10 coded diagnoses in MSDW collected since 2003
	370 million Number of test results , including: lab, radiology and pathology	189 million Number of medication records , including prescription and med admins	45 million Number of clinical documentations like progress notes, discharge summaries & operative reports

- Self service, de-identified access through the Cohort Query Tool (CQT), TriNetX, i2b2 Atlas/OMOP
- Identified data available by request (fee)
- Integrated with BioME and the Image Research Warehouse (>200,000 images)
- Creating a database of unstructured physician notes for Natural Language Processing (NLP) analysis (Clinithink with BioME)

A Recent Article You May Enjoy

BD04CH09_Percha ARjats.cls April 30, 2021 17:20

ANNUAL REVIEWS

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
The *Annual Review of Biomedical Data Science* is online at biodatasci.annualreviews.org
<https://doi.org/10.1146/annurev-biodatasci-030421-030931>
Copyright © 2021 by Annual Reviews.
All rights reserved.

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.

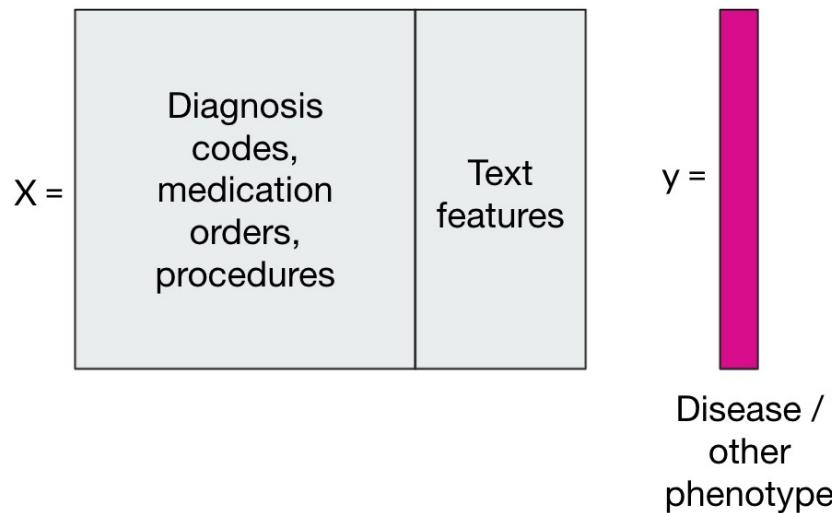
165

Part II

Defining the Task and Approach

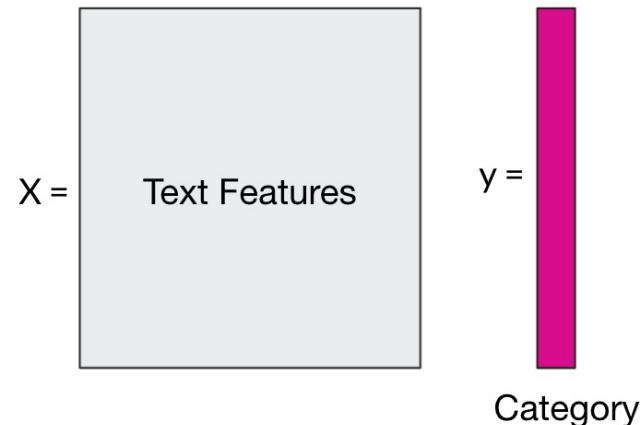
Information Extraction vs. Modeling

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. Some projects combine text data with other features, while others are about interpreting the text itself.



Example: Electronic phenotyping

EHR search indexing,
knowledgebase construction, patient
timeline building likewise focus on IE



Examples: Classify mammography reports by BI-RADS category, classify pathology reports by diagnosis (also many unsupervised examples)

Rule-Based vs. Statistical

Rule-based systems codify expert knowledge into a set of structured rules, or templates, that produce structured information when applied to unstructured text. Statistical (ML) methods can be difficult in the clinical domain because of privacy concerns and lack of training data.

Rule-based grammar	Machine Learning algorithm
 <ul style="list-style-type: none">★ Flexible★ Easy to debug★ Doesn't require a massive training corpus★ Understanding of the language phenomenon★ High precision	<ul style="list-style-type: none">★ Easy to scale★ "Learnability" without being explicitly programmed★ Fast development (if datasets available)★ High recall (coverage)
 <ul style="list-style-type: none">★ Requires skilled developers and linguists★ Slow parser development★ Moderate recall (coverage)	<ul style="list-style-type: none">★ Requires training corpus with annotation★ Difficult to debug★ No understanding of the language phenomenon

<https://medium.com/friendly-data/machine-learning-vs-rule-based-systems-in-nlp-5476de53c3b8>

Part III

Clinical Information Extraction

Three Common Tasks

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.

1. Named Entity Recognition (NER)

Identifying and locating mentions of conceptual categories, such as drug, symptom, or disease names, in text

2. 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. 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

Named Entity Recognition (NER)

The annotations shown here are from the Stanza clinical text processing pipeline, trained using data from the 2010 i2b2/VA challenge.

What do you notice?

Progress note:

Ms. S. is a 43F h/o antiphospholipid syndrome, HTN, DM, here for routine follow-up appointment.

Her main concern today is a 3 month h/o worsening shortness of breath a/w lower extremity edema. Her exercise tolerance has decreased from 10 blocks to half a block over the past few years.

She denies chest pain, palpitations, orthopnea, calf tenderness, fevers, and substance use.

On exam, she is afebrile, BP 110/80, HR 80, RR 18, O₂ Sat 98% on room air. She was AOx3, +JVD to her mid-neck, +HJ reflux, RRR, normal S1, prominent P2. Lungs CTAB, no wheezes, rales, or rhonchi. Abdomen soft, non-tender, non-distended, +bowel sounds. 1+ lower extremity pitting edema bilaterally to the shins.

Labs notable for Hgb 12, platelet count of 350, Na 135, K 3.5, Cr 1.

Chest X ray was clear.

EKG with a HR 84, NSR, no axis deviation, no ischemic changes.

Assessment:

Ms. S. is a 43F h/o antiphospholipid syndrome, HTN, DM, with progressive dyspnea and lower extremity edema concerning for new-onset acute decompensated heart failure. Given that she is afebrile without any infectious symptoms, pneumonia is less likely. She does not smoke which makes COPD exacerbation less likely.

Plan:

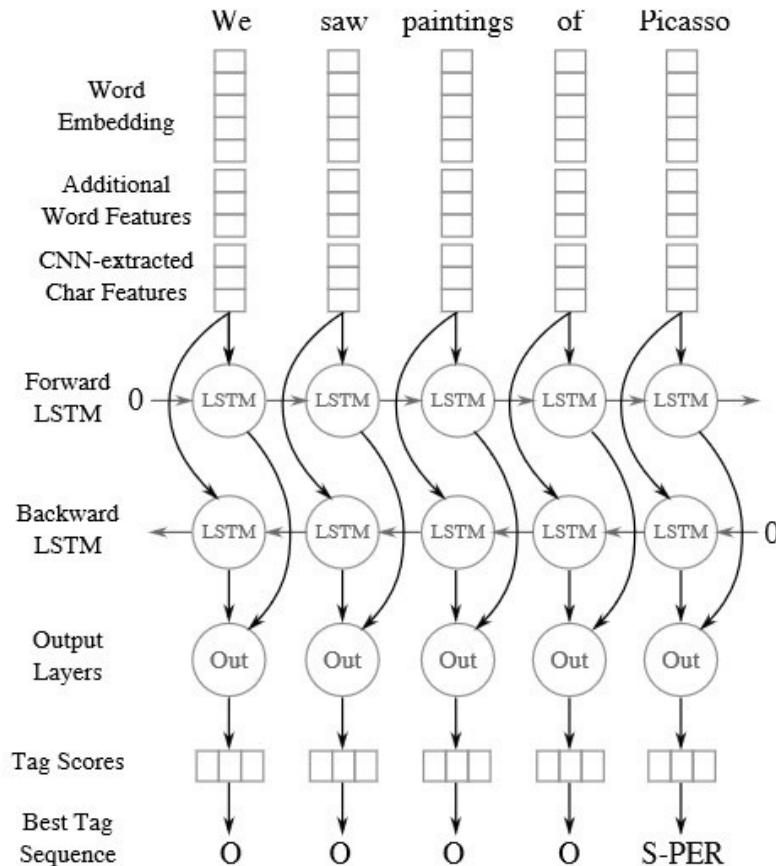
-BNP
-Transthoracic echocardiogram

Stanza clinical pipeline annotations

- Test
- Problem
- Treatment (none found)

Methods for NER

Systems range from simple dictionary-based approaches to 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).



<https://towardsdatascience.com/named-entity-recognition-ner-meeting-industrys-requirement-by-applying-state-of-the-art-deep-698d2b3b4ede>

Side Note...

Features for NER Algorithms

Tokenization

Tokenization is the act of splitting a text string into a list of tokens (e.g. words and punctuation, characters, etc.). The fastest tokenizers are based on regular expressions. For many applications, whitespace tokenization works just fine.

Original sentence:

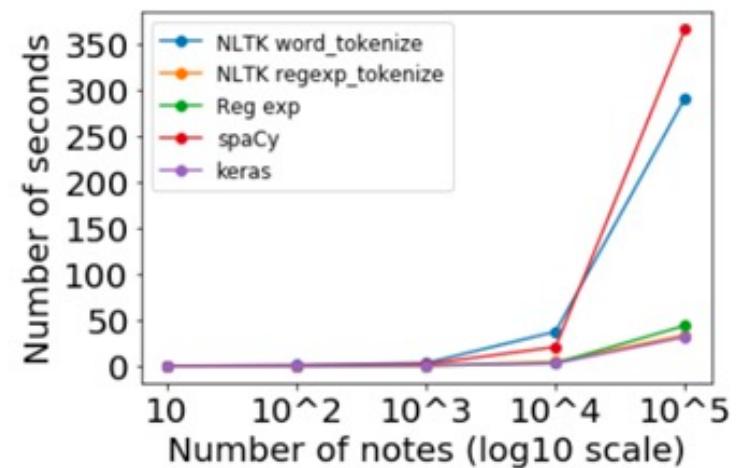
"Ethanol consumption reduces endothelium-dependent relaxation induced by acetylcholine."

Output of NLTK's word_tokenize:

"Ethanol", "consumption",
"reduces", "endothelium-dependent",
"relaxation", "induced", "by",
"acetylcholine", "."

Output of NLTK's wordpunct_tokenize:

"Ethanol", "consumption",
"reduces", "endothelium-dependent",
"relaxation", "induced", "by",
"acetylcholine", "."

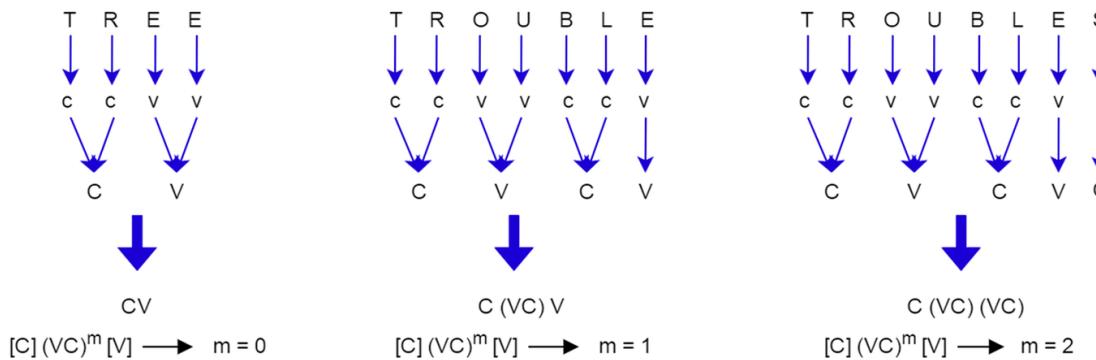


Source: <https://towardsdatascience.com/benchmarking-python-nlp-tokenizers-3ac4735100c5>

Stemming and Lemmatization

“The goal of both stemming and lemmatization is to reduce inflectional forms and sometimes derivationally related forms of a word to a common base form.”

– *Introduction to Information Retrieval (Manning and Schütze)* –



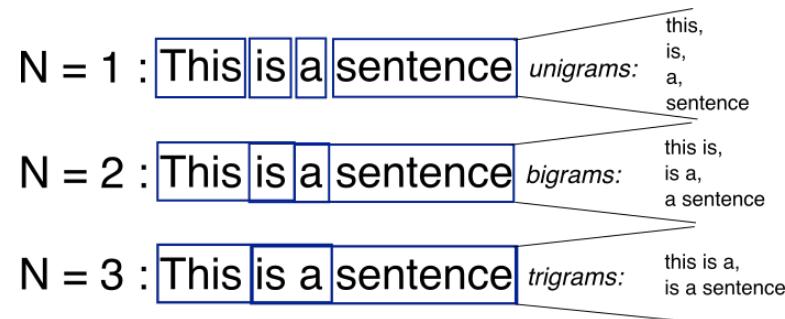
Porter Stemming Algorithm

SSES	→	SS	(m>0) ATIONAL	→	ATE
IES	→	I	(m>0) TIONAL	→	TION
SS	→	SS	(m>0) ENCI	→	ENCE
S	→		(m>0) ANCI	→	ANCE

Source: <https://vijinimallawaarachchi.com/2017/05/09/porter-stemming-algorithm/>

N-grams

A classic way of featurizing documents: collect all words and phrases of length N=1, 2, 3, etc. The combinatorics typically become a problem above N=3 (sometimes even N=2). N-grams can also be produced at the character level.



- unigram (1-gram):

a	swimmer	likes	swimming	thus	he	swims
---	---------	-------	----------	------	----	-------

- bigram (2-gram):

a swimmer	swimmer likes	likes swimming	swimming thus	...
-----------	---------------	----------------	---------------	-----

- trigram (3-gram):

a swimmer likes	swimmer likes swimming	likes swimming thus	...
-----------------	------------------------	---------------------	-----

TFIDF

How should one select the top features for downstream use in classification algorithms? It's not good to use pure frequency because the most frequent N-grams are often the least informative.

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

$tf_{i,j}$ = number of occurrences of i in j

df_i = number of documents containing i

N = total number of documents

Sentence 1: The car is driven on the road.

Sentence 2: The truck is driven on the highway.

Word	TF		IDF	TF*IDF	
	A	B		A	B
The	1/7	1/7	$\log(2/2) = 0$	0	0
Car	1/7	0	$\log(2/1) = 0.3$	0.043	0
Truck	0	1/7	$\log(2/1) = 0.3$	0	0.043
Is	1/7	1/7	$\log(2/2) = 0$	0	0
Driven	1/7	1/7	$\log(2/2) = 0$	0	0
On	1/7	1/7	$\log(2/2) = 0$	0	0
The	1/7	1/7	$\log(2/2) = 0$	0	0
Road	1/7	0	$\log(2/1) = 0.3$	0.043	0
Highway	0	1/7	$\log(2/1) = 0.3$	0	0.043

Source: <https://www.freecodecamp.org/news/how-to-process-textual-data-using-tf-idf-in-python-cd2bbc0a94a3/>

Feature Selection and Normalization

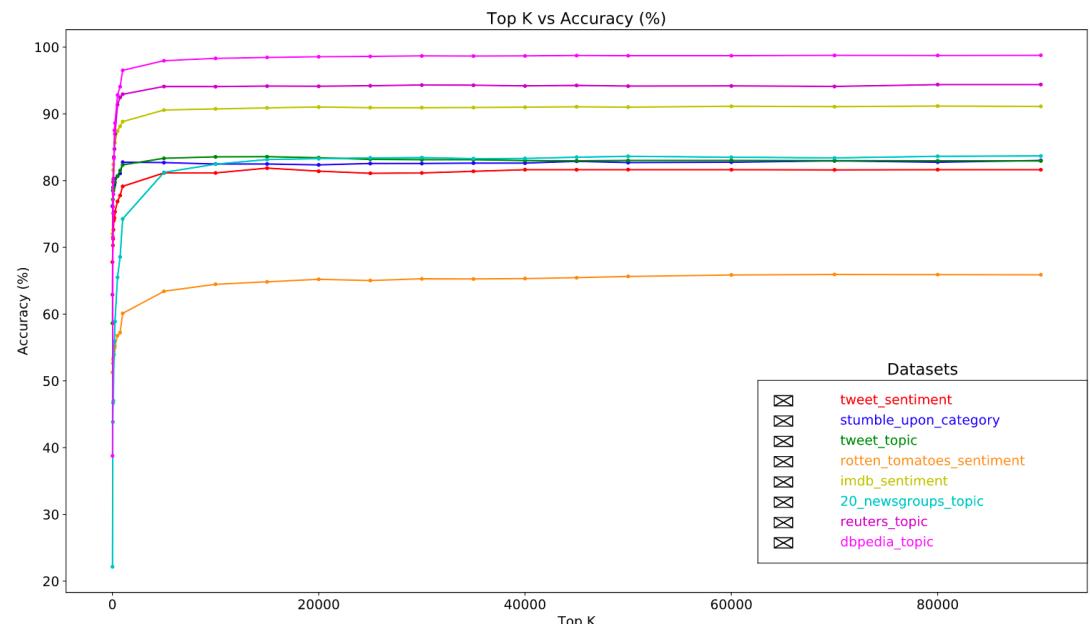
Text classification accuracy tends not to benefit much from features beyond the top ~20k. You may choose to **normalize** features by sample or by feature or do nothing (most common).

sklearn.feature_selection: Feature Selection

The `sklearn.feature_selection` module implements feature selection algorithms. It currently includes univariate filter selection methods and the recursive feature elimination algorithm.

User guide: See the [Feature selection](#) section for further details.

<code>feature_selection.GenericUnivariateSelect(...)</code>	Univariate feature selector with configurable strategy.
<code>feature_selection.SelectPercentile(...)</code>	Select features according to a percentile of the highest scores
<code>feature_selection.SelectKBest([score_func, k])</code>	Select features a
<code>feature_selection.SelectFpr([score_func, alpha])</code>	Filter: Select the
<code>feature_selection.SelectFdr([score_func, alpha])</code>	Filter: Select the
<code>feature_selection.SelectFromModel(estimator)</code>	Meta-transformer weights.
<code>feature_selection.SelectFwe([score_func, alpha])</code>	Filter: Select the
<code>feature_selection.RFE(estimator[, step, ...])</code>	Feature ranking \n validated selectio
<code>feature_selection.RFECV(estimator[, step, ...])</code>	Feature ranking \n validated selectio
<code>feature_selection.VarianceThreshold([threshold])</code>	Feature selector
<code>feature_selection.chi2(X, y)</code>	Compute chi-squa and class.
<code>feature_selection.f_classif(X, y)</code>	Compute the ANC
<code>feature_selection.f_regression(X, y[, center])</code>	Univariate linear re
<code>feature_selection.mutual_info_classif(X, y)</code>	Estimate mutual in
<code>feature_selection.mutual_info_regression(X, y)</code>	Estimate mutual in

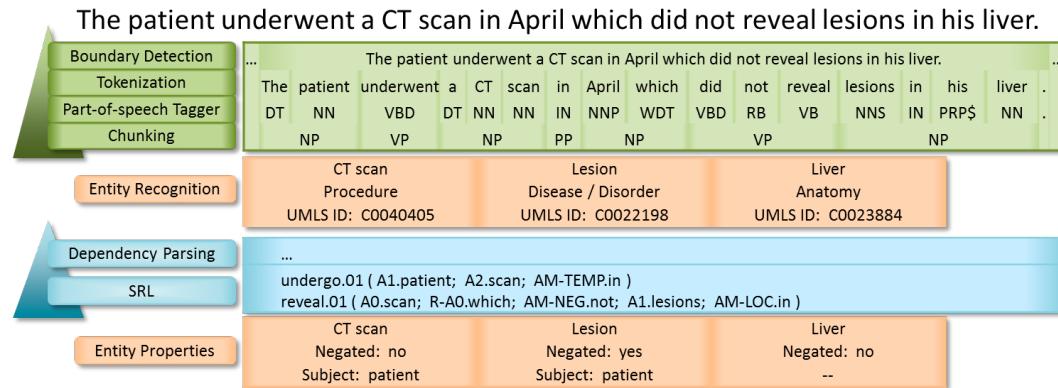


Source: <https://developers.google.com/machine-learning/guides/text-classification/step-3>

Biomedical Feature Extraction and UMLS

There are a variety of biomedical text annotators available that will extract biomedically or clinically relevant concepts from free text. Most of these map text strings to concepts from the Unified Medical Language System (UMLS). **They do not always improve classification performance.**

cTAKES
<https://ctakes.apache.org/>



MetaMapLite
<https://metamap.nlm.nih.gov/MetaMapLite.shtml>

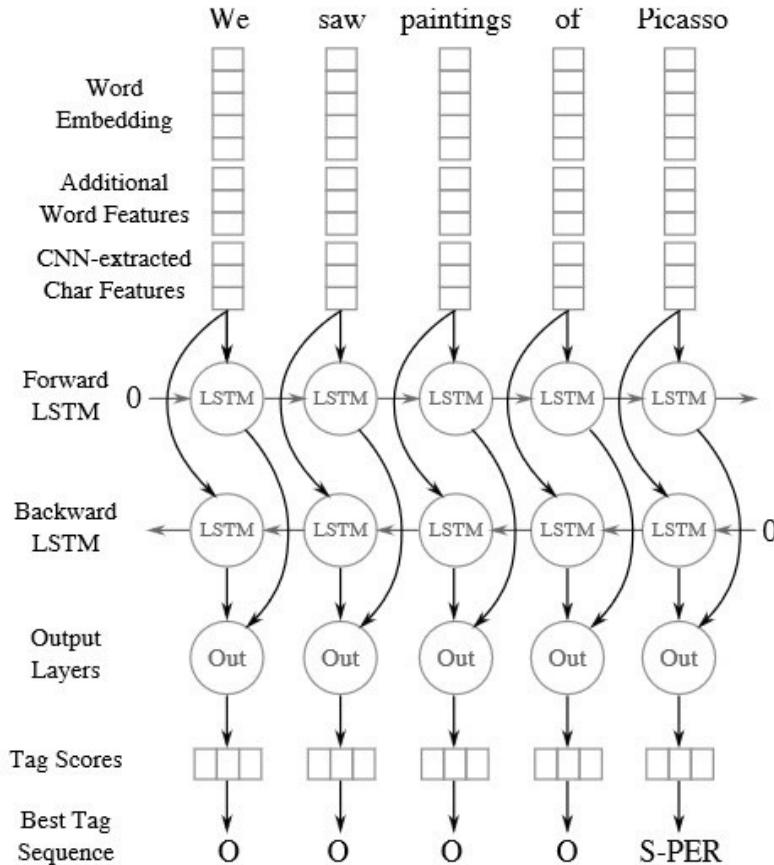
Sentences	Detected phrases	Extracted concepts	Selected
First sentence	Hyperlipidaemia	[Disease or Syndrome] [Finding]	✓
	Patient	[Patient or Disabled group]	✗
	Lipitor	[Organic Chemical, Pharmacologic Substance]	✗
	80%	[Quantitative Concept]	✗
	mg++ increased	[Finding]	✗
Second sentence	Progress note	[Clinical Attribute] [Intellectual Product]	✗ ✗
	Patient chart	[Manufactured Object]	✗
	Assisted living facility	[Healthcare Related Organization, Manufactured Object]	✗
	Patient	[Patient or Disabled group]	✗
	Shortness of breath	[Sign or Symptom]	✓

Source: Abdollahi et al. (2018). Uncovering Discriminative Knowledge-Guided Medical Concepts for Classifying Coronary Artery Disease Notes: 31st Australasian Joint Conference, Wellington, New Zealand, December 11-14, 2018.

Back to Part III...

Methods for NER (review)

Systems range from simple dictionary-based approaches to 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).



<https://towardsdatascience.com/named-entity-recognition-ner-meeting-industrys-requirement-by-applying-state-of-the-art-deep-698d2b3b4ede>

Domain Specificity and Key Challenges

NER only makes sense when the entities involved are discrete and have defined locations in text.

Not all NER annotations are correct or meaningful without consideration of the surrounding context.

One may be interested in an entity class for which no pre-annotated corpus or pretrained model is available.

General domain NER models trained using general domain corpora – may not work on clinical text.

Progress note:

Ms. S. is a 43F h/o antiphospholipid syndrome, HTN, DM, here for routine follow-up appointment.

Her main concern today is a 3 month h/o worsening shortness of breath a/w lower extremity edema. Her exercise tolerance has decreased from 10 blocks to half a block over the past few years.

She denies chest pain, palpitations, orthopnea, calf tenderness, fevers, and substance use.

On exam, she is afebrile, BP 110/80, HR 80, RR 18, O₂ Sat 98% on room air. She was AOx3, +JVD to her mid-neck, +HJ reflux, RRR, normal S1, prominent P2. Lungs CTAB, no wheezes, rales, or rhonchi. Abdomen soft, non-tender, non-distended, +bowel sounds. 1+ lower extremity pitting edema bilaterally to the shins.

Labs notable for Hgb 12, platelet count of 350, Na 135, K 3.5, Cr 1.

Chest X ray was clear.

EKG with a HR 84, NSR, no axis deviation, no ischemic changes.

Assessment:

Ms. S. is a 43F h/o antiphospholipid syndrome, HTN, DM, with progressive dyspnea and lower extremity edema concerning for new-onset acute decompensated heart failure. Given that she is afebrile without any infectious symptoms, pneumonia is less likely. She does not smoke which makes COPD exacerbation less likely.

Plan:

-BNP
-Transthoracic echocardiogram

Stanza clinical pipeline annotations

- Test
- Problem
- Treatment (none found)

Concept Normalization

Also known as “entity linking”, this is the task of assigning a unique identity to each entity name mentioned in a text.

What do you notice?

Progress note:

Ms. S. is a 43F h/o antiphospholipid syndrome HTN DM, here for routine follow-up appointment.

Her main concern today is a 3 month h/o worsening shortness of breath a/w lower extremity edema. Her exercise tolerance has decreased from 10 blocks to half a block over the past few years.

She denies chest pain, palpitations, orthopnea, calf tenderness, fevers, and substance use.

On exam, she is afebrile, BP 110/80, HR 80, RR 18, O2 Sat 98% on room air. She was Aox3, +JVD to her mid-neck, +HJ reflux RRR, normal S1, prominent P2. Lungs CTAB, no wheezes, rales, or rhonchi. Abdomen soft, non-tender, non-distended, +bowel sounds.

cTAKES default pipeline annotations

- MedicationMention
- DiseaseDisorderMention
- SignSymptomMention
- ProcedureMention
- AnatomicalSiteMention
- Negated
- Uncertainty detected

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	Gastresophageal reflux disease	C0017168
11	MedicationMention	CTAB	Cetrimonium bromide	C0951233
23	MedicationMention	BNP	Nesiritide	C0054015

Plan:

-BNP

-Transthoracic echocardiogram

Numbers, Ranges, and Sections

Identifying and normalizing numeric values and identifying note sections are tasks somewhat unique to clinical text. To date, the SecTag system developed by Denny et al is the only section identification system used outside the institution in which it was developed.

Research article | [Open Access](#) | Published: 18 July 2019

Current approaches to identify sections within clinical narratives from electronic health records: a systematic review

[Alexandra Pomares-Quimbaya](#) , [Markus Kreuzthaler](#) & [Stefan Schulz](#)

[BMC Medical Research Methodology](#) **19**, Article number: 155 (2019) | [Cite this article](#)

3864 Accesses | 2 Citations | 1 Altmetric | [Metrics](#)

Abstract

Background

The identification of sections in narrative content of Electronic Health Records (EHR) has demonstrated to improve the performance of clinical extraction tasks; however, there is not yet a shared understanding of the concept and its existing methods. The objective is to report the results of a systematic review concerning approaches aimed at identifying sections in narrative content of EHR, using both automatic or semi-automatic methods.

Methods

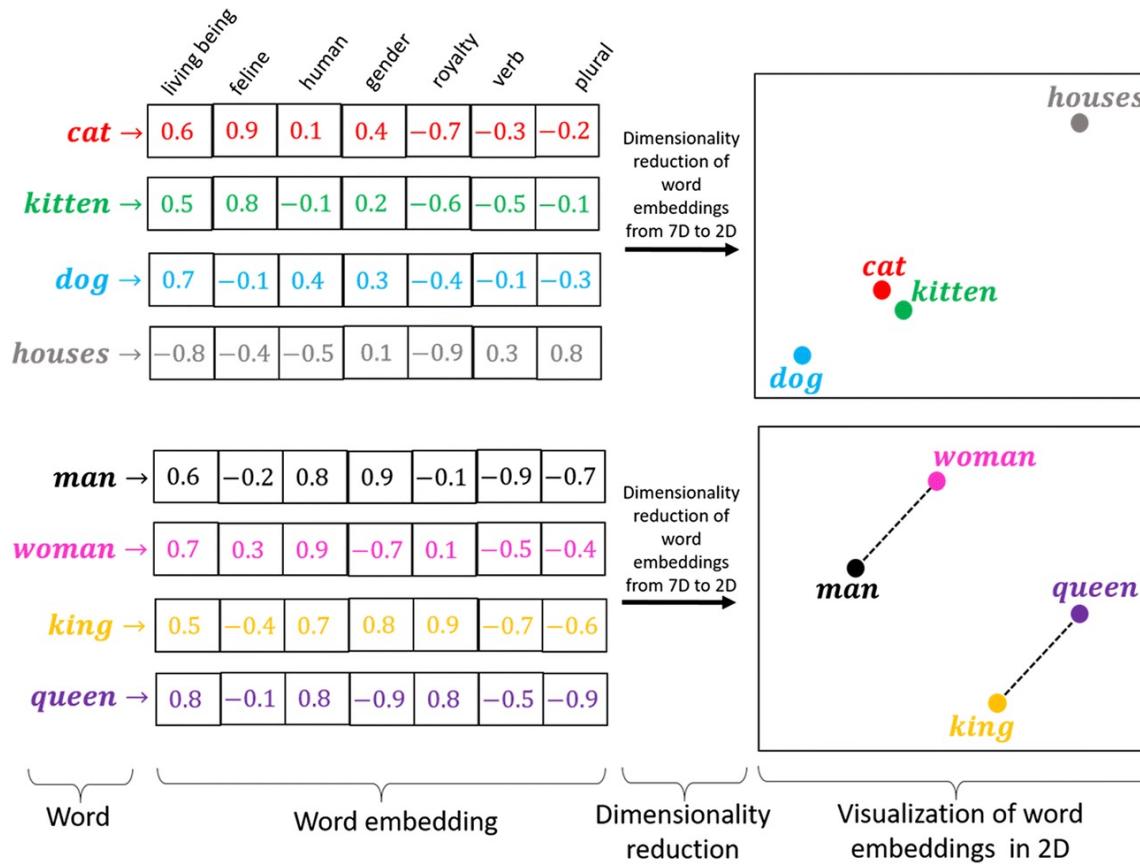
This review includes articles from the databases: SCOPUS, Web of Science and PubMed (from January 1994 to September 2018). The selection of studies was done using predefined

Part IV

Embeddings and Pretraining

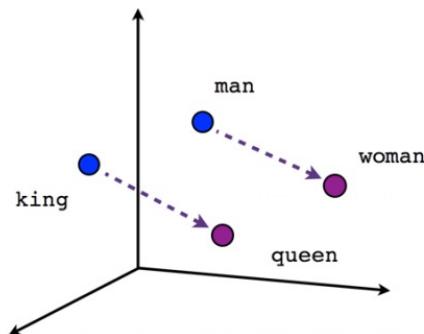
Embeddings

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.

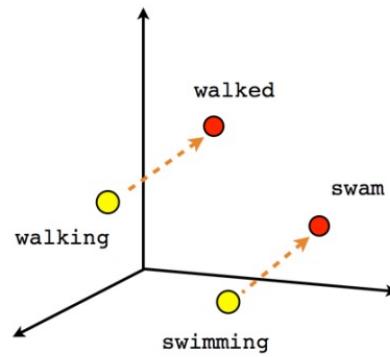


Embeddings are not a New Idea

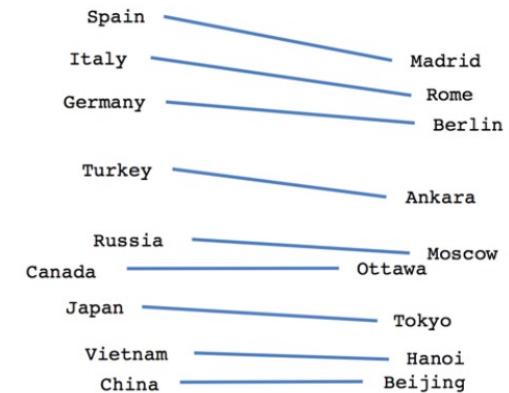
The field of distributional semantics originated with Latent Semantic Indexing in the late 1980s and reached a milestone with the development of word2vec and GloVe in 2013-2014. It is a collection of methods that create vector space embeddings of words and phrases that reflect how they are used in context.



Male-Female



Verb tense



Country-Capital

Why are embeddings useful?

(1) Do not require annotated corpora for training, so it is easy to create embeddings that are specific to clinical text or that capture regularities of expression within a clinical subfield or institution. (2) An embedding can incorporate structured information beyond what is found in text, and embeddings have been created to represent CUIs, documents, or entire patient records.

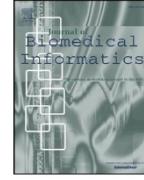
Journal of Biomedical Informatics 101 (2020) 103323

Contents lists available at ScienceDirect

Journal of Biomedical Informatics

journal homepage: www.elsevier.com/locate/yjbin

 ELSEVIER



SECNLP: A survey of embeddings in clinical natural language processing

Katikapalli Subramanyam Kalyan*, S. Sangeetha



Text Analytics and NLP Lab, Department of Computer Applications, NIT Trichy, India

ARTICLE INFO

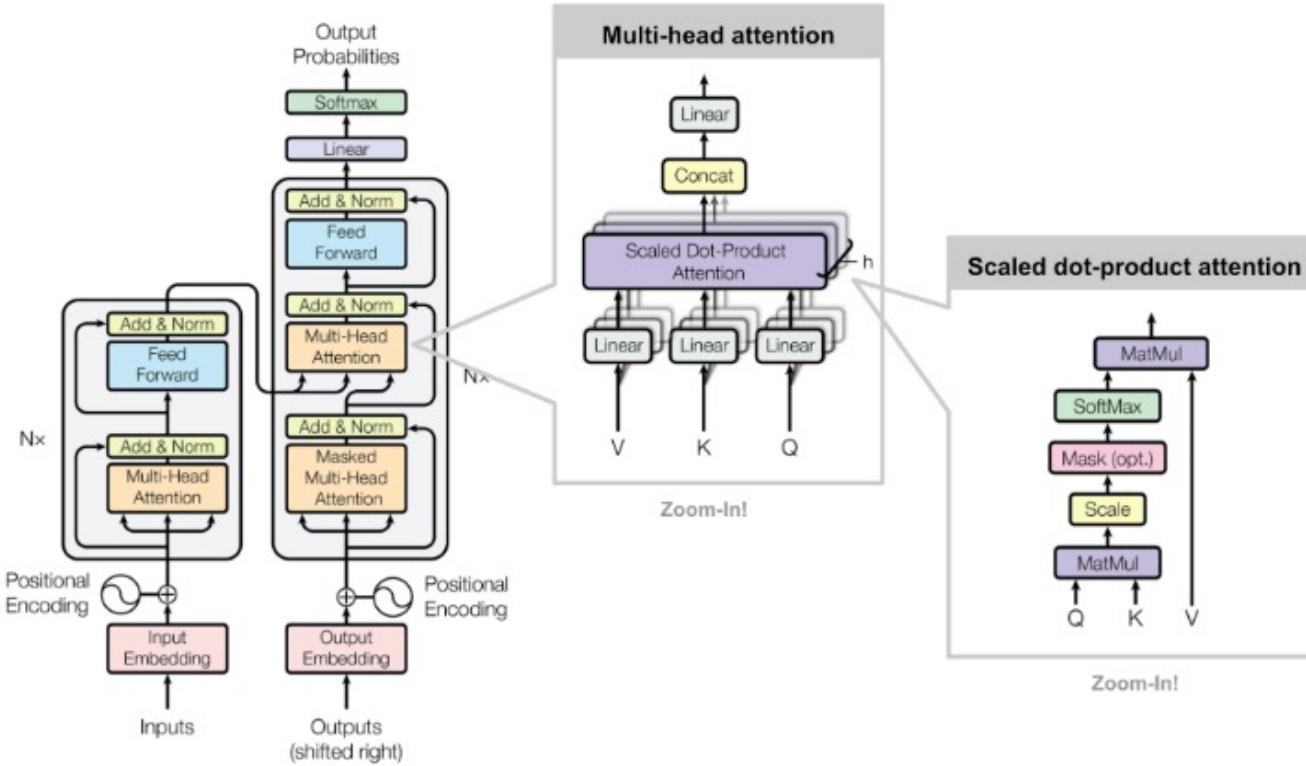
Keywords:
Embeddings
Distributed representations
Medical
Natural language processing
Survey

ABSTRACT

Distributed vector representations or embeddings map variable length text to dense fixed length vectors as well as capture prior knowledge which can be transferred to downstream tasks. Even though embeddings have become de facto standard for text representation in deep learning based NLP tasks in both general and clinical domains, there is no survey paper which presents a detailed review of embeddings in Clinical Natural Language Processing. In this survey paper, we discuss various medical corpora and their characteristics, medical codes and present a brief overview as well as comparison of popular embeddings models. We classify clinical embeddings and discuss each embedding type in detail. We discuss various evaluation methods followed by possible solutions to various challenges in clinical embeddings. Finally, we conclude with some of the future directions which will advance research in clinical embeddings.

Contextual Embeddings and Pretraining

Novel neural network architectures have permitted the creation of embeddings that vary depending on the context. This has expanded the representational power of embedding methods.



<https://medium.com/@shreyasikalra25/predict-movie-reviews-with-bert-88d8b79f5718>

Transformer Architectures

The best-known Transformer-based model is called BERT. It was developed by a team at Google and published in 2019.



ULMfit	GPT	BERT	GPT-2
Jan 2018	June 2018	Oct 2018	Feb 2019
Training: 1 GPU day	Training 240 GPU days	Training 256 TPU days ~320–560 GPU days	Training ~2048 TPU v3 days according to a reddit thread



Source: Course notes for CS224N (Stanford). Lecture 14: Contextual Representations and Pretraining

A Good Place to Start

If you really want to understand neural network methods for natural language processing, read this paper first (and then study all the transformer-based models that have come out in the last three years).

Journal of Artificial Intelligence Research 57 (2016) 345–420

Submitted 9/15; published 11/16

A Primer on Neural Network Models for Natural Language Processing

Yoav Goldberg

*Computer Science Department
Bar-Ilan University, Israel*

YOAV.GOLDBERG@GMAIL.COM

Abstract

Over the past few years, neural networks have re-emerged as powerful machine-learning models, yielding state-of-the-art results in fields such as image recognition and speech processing. More recently, neural network models started to be applied also to textual natural language signals, again with very promising results. This tutorial surveys neural network models from the perspective of natural language processing research, in an attempt to bring natural-language researchers up to speed with the neural techniques. The tutorial covers input encoding for natural language tasks, feed-forward networks, convolutional networks, recurrent networks and recursive networks, as well as the computation graph abstraction for automatic gradient computation.

1. Introduction

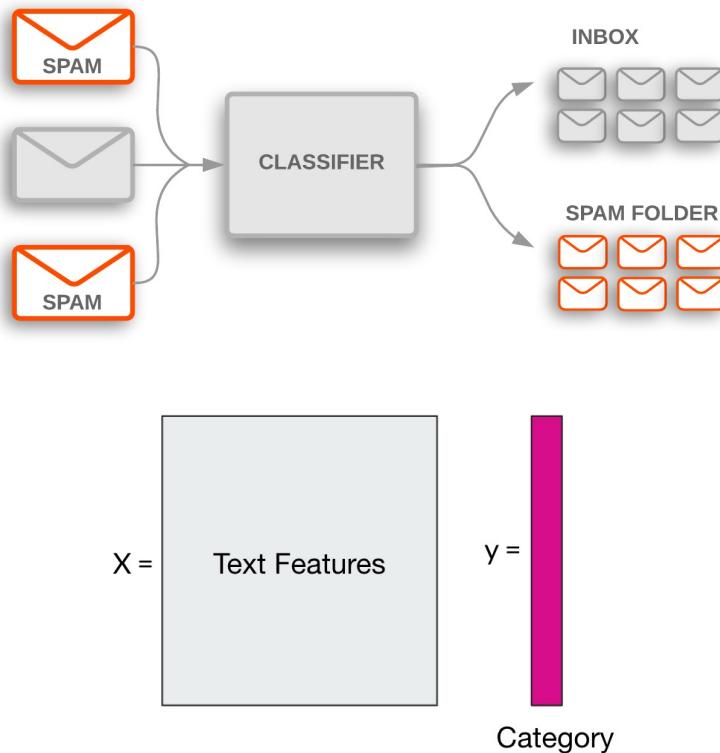
For over a decade, core NLP techniques were dominated by machine-learning approaches that used linear models such as support vector machines or logistic regression, trained over very high dimensional yet very sparse feature vectors.

Part V

Text Classification

Text Classification

The goal is to classify the text of a document (where a document can be as short as a single phrase or sentence and as long as a book) into two or more categories. A recent survey found that of 212 recent clinical text mining papers, 88 (41.5%) focused on text classification.



Examples:

- Classifying primary care descriptions of back pain into acute vs. chronic
- Distinguishing normal vs. abnormal knee MRI reports
- Assessing whether a patient is a current or former smoker based on clinical notes

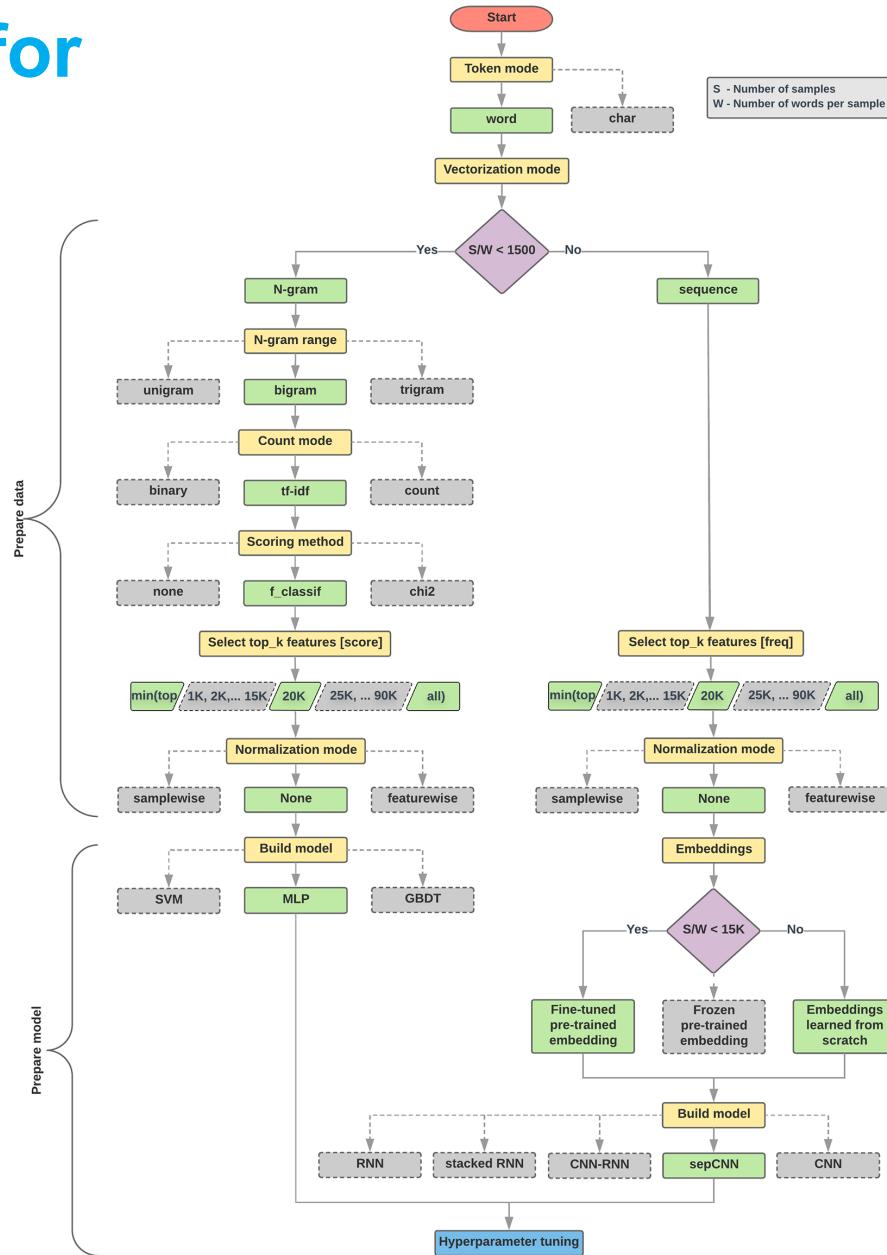
General Procedure for Text Classification

Any machine learning classifier can be adapted to classify text.

The hard part is figuring out how to turn a document into a set of features suitable for classification.

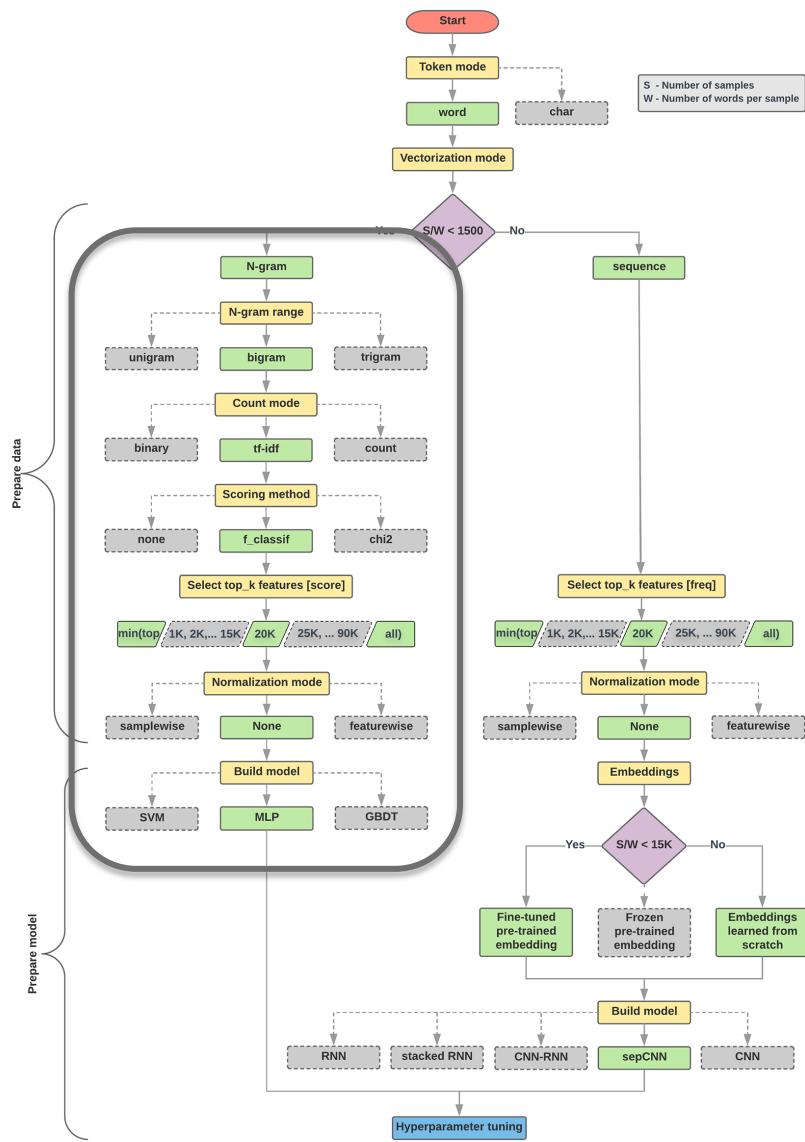
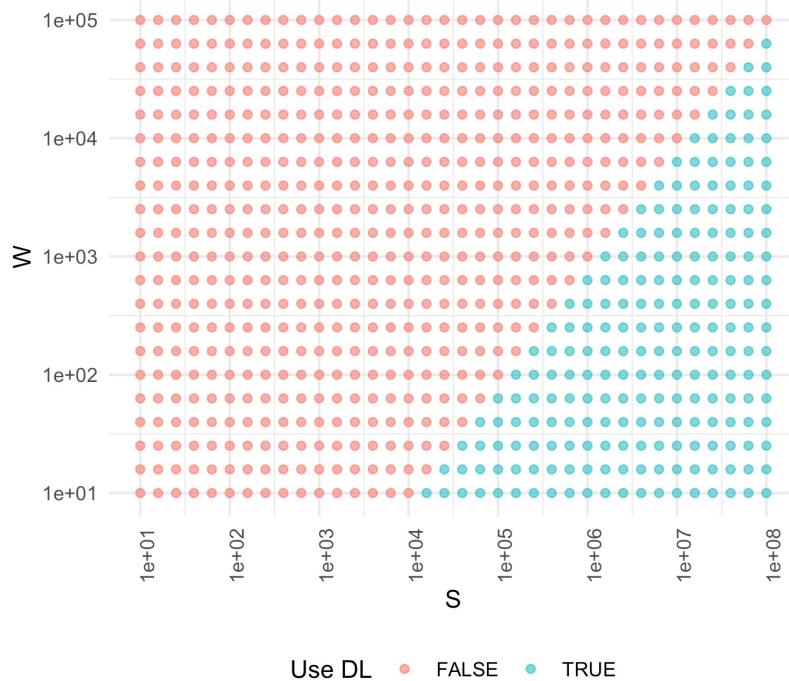
Google provides this handy guide →

<https://developers.google.com/machine-learning/guides/text-classification>



The Standard Pipeline

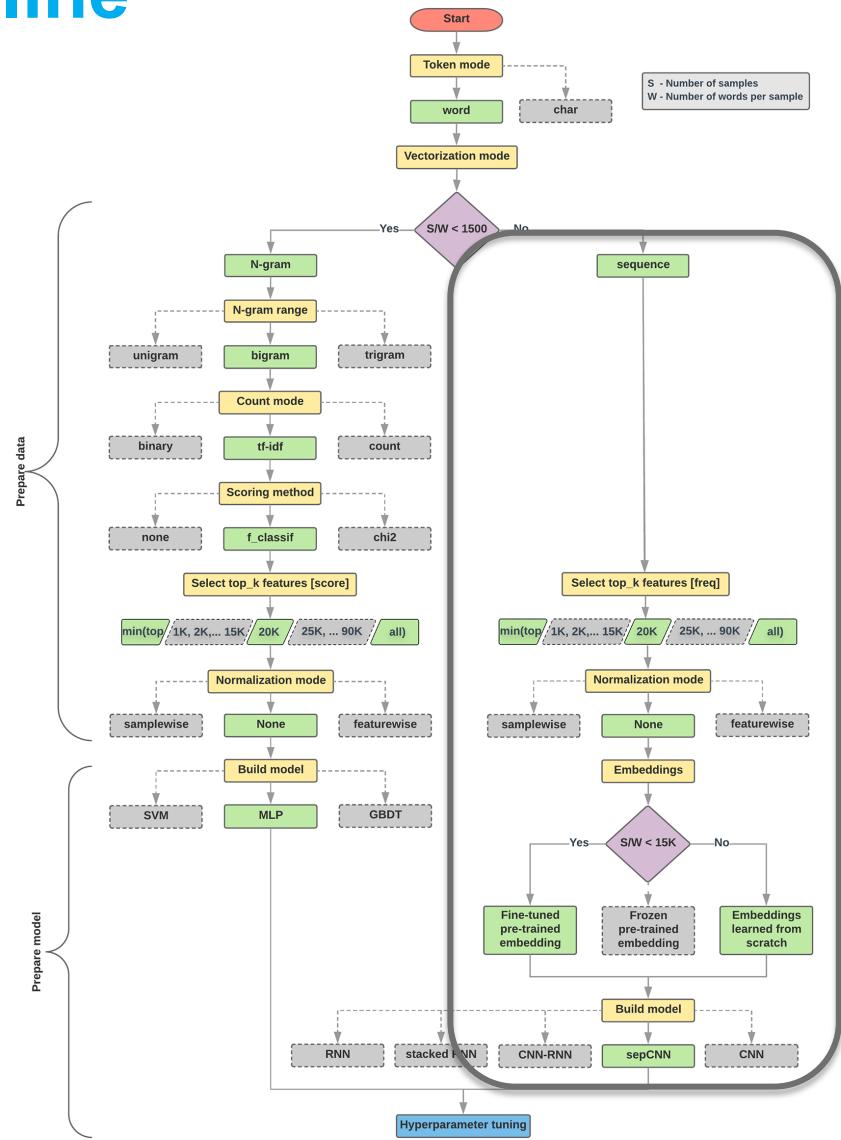
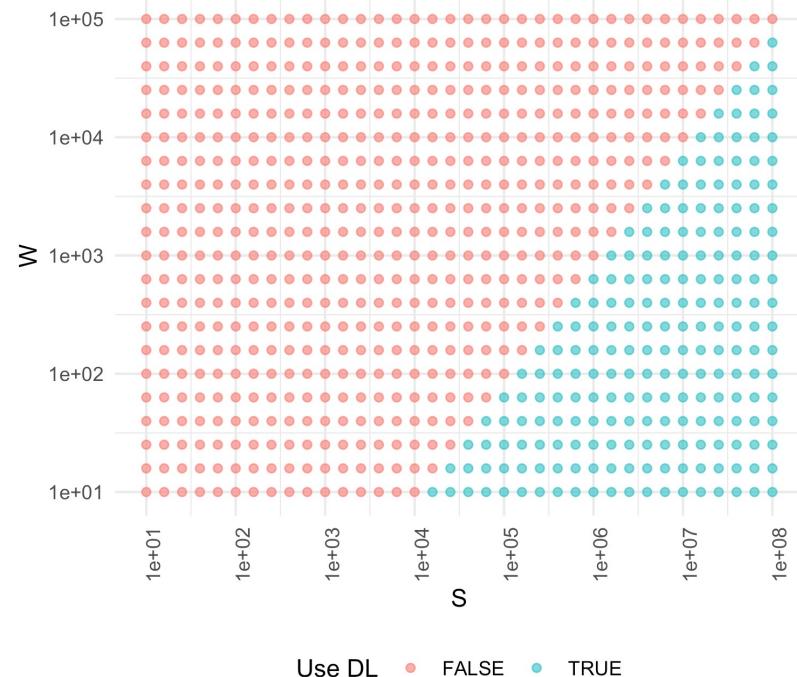
Start with raw text, tokenize it, extract N-grams, weight them somehow (or not), select top features, normalize them (or not), use whatever supervised ML algorithm you want to classify them.



The Deep Learning Pipeline

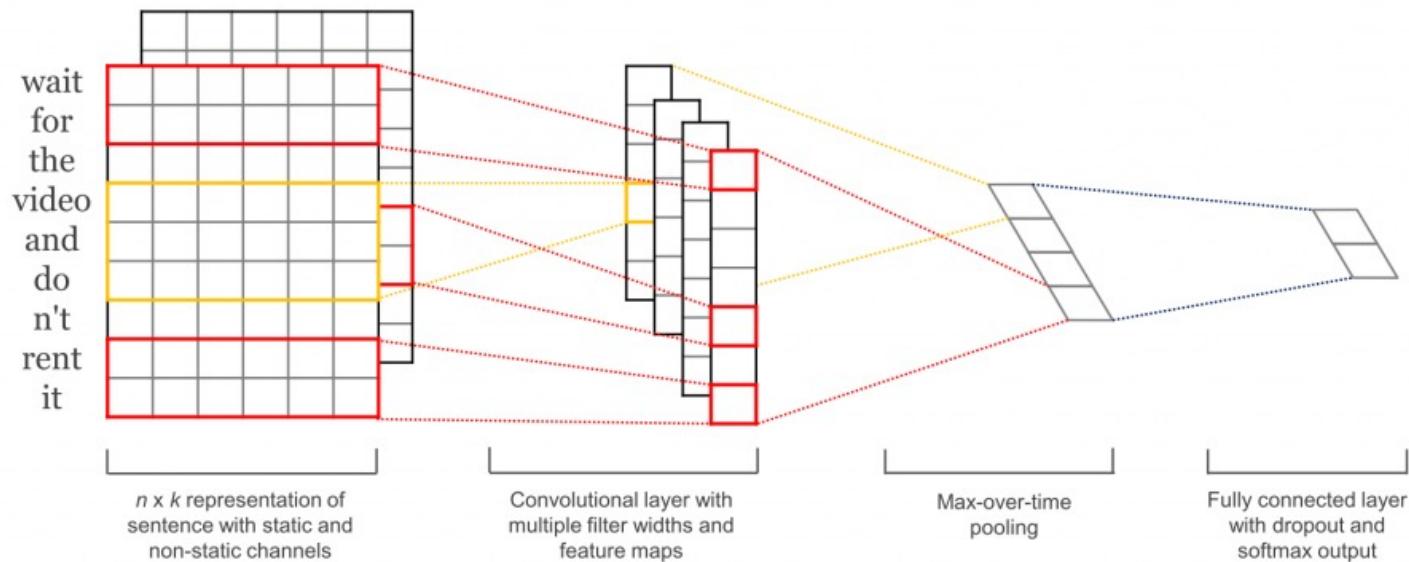
Start with raw text, tokenize it, do feature selection (maybe), produce embeddings or use pretrained embeddings, feed to neural network.

Convolutional neural networks (CNNs) are most popular architecture at the moment.



Text Classification with CNNs

These models have been deployed on a variety of tasks, including assigning diagnosis codes, classifying radiology reports, subtyping diseases, and determining the presence or absence of comorbidities.



<http://www.wildml.com/2015/12/implementing-a-cnn-for-text-classification-in-tensorflow/>

Part VI

Additional Topics: Weak Supervision
and Inference/Entailment

Weak and Distant Supervision

Machine learning approaches to clinical NLP generally suffer from a lack of training data. Existing clinical information extraction and text classification models have generally been trained using the same few annotated datasets, which restricts the range and quality of annotations they produce.

1. 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

2. 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 data associated with patients/encounters, not sentences/documents)

3. Alternatives: Crowdsourcing and Active Learning

Crowdsourcing not usually a viable option due to privacy concerns.
Active learning is a strategy for minimizing annotation effort by iteratively sampling subsets of data for human annotation based on current predictions of a supervised learning algorithm.

Relation Extraction

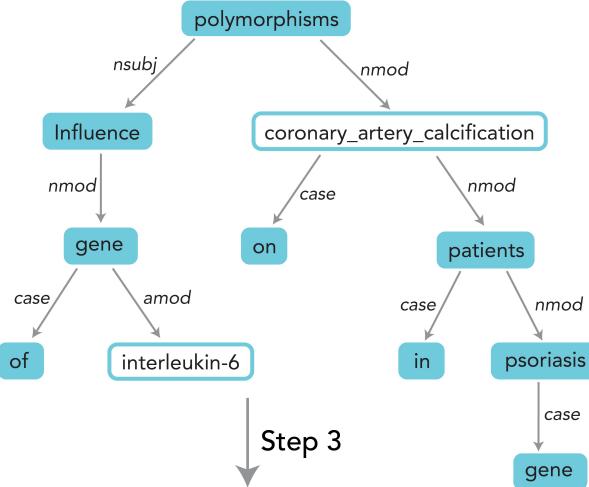
Relation extraction is the task of assigning a structured form to a relationship between or among entities based on how they are described in text. Typically, this form includes the categories of the entities and a label denoting the nature of their relationship.

Influence of interleukin-6 gene polymorphisms on coronary_artery_calcification in patients with psoriasis.

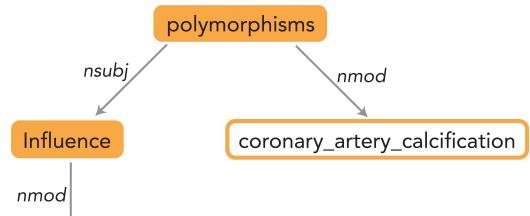
Step 1

Influence of interleukin-6 gene polymorphisms on coronary_artery_calcification in patients with psoriasis.

Step 2



Step 3



[interleukin-6, amod, gene, nmod, influence, nsubj, polymorphisms, amod, coronary_artery_calcification]

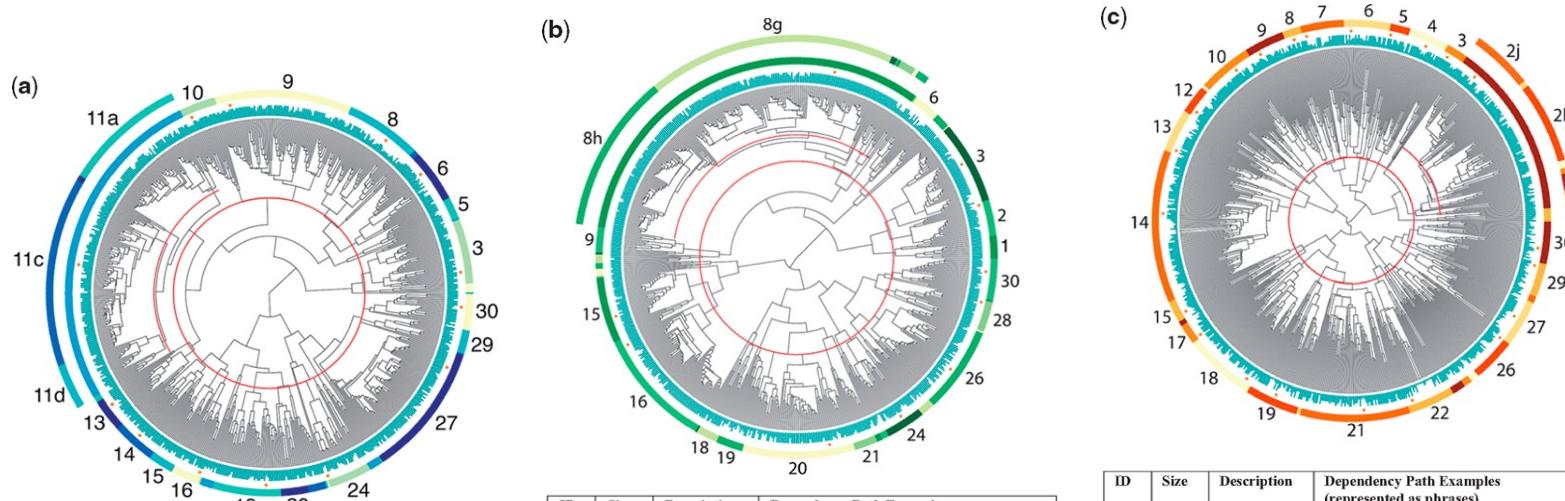
Step 4



Percha B, Altman RB (2018) A global network of biomedical relationships derived from text. *Bioinformatics* 34(15), 2614-2624. Figure 1.

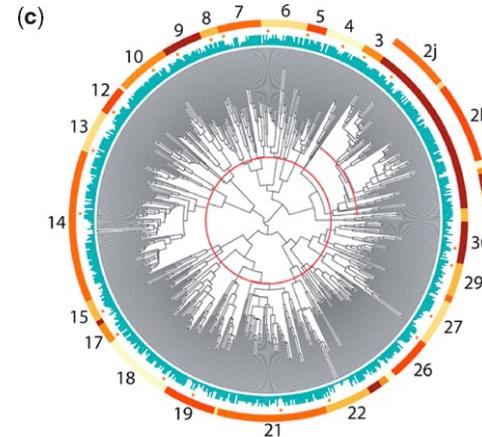
Example from the Biomedical Domain: GNBR

Cluster patterns from text (dependency paths) into semantically meaningful categories.



ID	Size	Description	Dependency Path Examples (represented as phrases)
3	36	inhibition	<ul style="list-style-type: none"> "C, a G inhibitor" "G specific inhibitor, C" "C, an inhibitor of G" "G inhibition by C" "effects of the G inhibitor, C, on..."
11c	96	metabolism, secretion/uptake	<ul style="list-style-type: none"> "effect of G on C metabolism" "effects of G on C formation" "effect of G on the secretion of C" "control of G by C" "G stimulates C uptake" "[chemical] may reduce G concentration via C" "G stimulates C transport" "effect of C on G release"
19	38	channels / transporters	<ul style="list-style-type: none"> "regulation of G transporters C and..." "G is a C channel that modulates..." "G C transporter expression" "C transporter, G" "distribution of G C channel subunits"

ID	Size	Description	Dependency Path Examples (represented as phrases)
8h	80	treatment of disease (indication of efficacy)	<ul style="list-style-type: none"> "C may be useful for the treatment of D" "evaluate the protective efficacy of C in D" "C is a promising treatment option for patients with D" "C is approved for the treatment of D" "C is commonly prescribed for D"
15	37	side effects (association)	<ul style="list-style-type: none"> "D associated with C therapy" "the use of C has been associated with D" "C intake was associated with D" "incidence of D in patients receiving C" "D occurred after C"
20	63	levels associated with disease risk / progression	<ul style="list-style-type: none"> "high C levels are associated with increased risk of D" "C implicated in D" "effect of D on serum C levels" "patients with D and increased C concentrations" "C has been implicated in the pathogenesis of D" "C intake may be associated with [lower/higher] risk of D" "C supplementation and incidence of D: ..."



ID	Size	Description	Dependency Path Examples (represented as phrases)
2j	33	influences disease treatment	<ul style="list-style-type: none"> "the use of G in the treatment of D" "D in patients treated with G" "effect of G on [event] in D patients" "G therapy in patients with D" "efficacy of G in D"
7	26	biomarkers, diagnostic	<ul style="list-style-type: none"> "G is a robust diagnostic biomarker for D" "G is an independent predictor of D" "G as an indicator of D in patients with..." "prognostic significance of G in D patients" "effects of [situation/event] on G levels in D" "G is a potential marker of D"
14	91	causal mutations	<ul style="list-style-type: none"> "mutation of G in a patient with D" "G mutation is associated with D" "novel mutation in G gene associated with D" "characterization of G mutations causing D" "mutations of the G gene in patients with D" "D: a novel G mutation..." "D: novel G mutations and..." "the recurrent mutation of G in C patients" "G mutations can cause D"

Percha B, Altman RB (2018) A global network of biomedical relationships derived from text. *Bioinformatics* 34(15), 2614-2624. Figure 3.

Why Relation Extraction is Important

Addresses issue of 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.

Progress note:

Ms. S. is a 43F h/o antiphospholipid syndrome HTN DM, here for routine follow-up appointment.

Her main concern today is a 3 month h/o worsening shortness of breath a/w lower extremity edema. Her exercise tolerance has decreased from 10 blocks to half a block over the past few years.

She denies chest pain, palpitations, orthopnea, calf tenderness, fevers, and substance use.

On exam, she is afebrile, BP 110/80, HR 80, RR 18, O2 Sat 98% on room air. She was Aox3, +JVD to her mid-neck, +HJ reflux RRR, normal S1, prominent P2. Lungs CTAB, no wheezes, rales, or rhonchi. Abdomen soft, non-tender, non-distended, +bowel sounds. 1+ lower extremity pitting edema bilaterally to the shins.

Labs notable for Hgb 12, platelet count of 350. Na 135, K 3.5, Cr 1.

Chest X ray was clear.

EKG with a HR 84, NSR, no axis deviation, no ischemic changes.

Assessment:

Ms. S. is a 43F h/o antiphospholipid syndrome HTN DM, with progressive dyspnea and lower extremity edema concerning for new-onset acute decompensated heart failure. Given that she is afebrile without any infectious symptoms, pneumonia is less likely. She does not smoke which makes COPD exacerbation less likely.

Plan:

-BNP
-Transthoracic echocardiogram

cTAKES default pipeline annotations

MedicationMention
DiseaseDisorderMention
SignSymptomMention
ProcedureMention
AnatomicalSiteMention
<input checked="" type="checkbox"/> Negated
<input type="checkbox"/> Uncertainty detected

Some Opinionated Commentary

- State-of-the-art machine learning models, which are very data hungry, are likely to hit a dead end when it comes to relation extraction and inference.
- Transfer learning with massive language models is a promising approach, but it is still basically pattern matching.
- Text mining is the key to unlocking the potential of electronic health records, esp. for generating “real world evidence”.
- There is a lot of money to be made if someone can crack this. Health technology companies often imply they are using NLP for data curation, but this is mostly a lie.

Great Reviews of the Field

A REVIEW OF REVIEWS

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

Year	Author(s)	Title	Reference
2011	Chapman <i>et al</i>	Overcoming Barriers to NLP for Clinical Text: The Role of Shared Tasks and the Need for Additional Creative Solutions	(141)
2016	Ford <i>et al</i>	Extracting Information from the Text of Electronic Medical Records to Improve Case Detection: A Systematic Review	(142)
2016	Koleck <i>et al</i>	Natural Language Processing of Symptoms Documented in Free-Text Narratives of Electronic Health Records: A Systematic Review	(5)
2017	Kreimeyer <i>et al</i>	Natural Language Processing Systems for Capturing and Standardizing Unstructured Clinical Information: A Systematic Review	(8)
2019	Khattak <i>et al</i>	A Survey of Word Embeddings for Clinical Text	(143)
2019	Mujtaba <i>et al</i>	Clinical Text Classification Research Trends: Systematic Literature Review and Open Issues	(18)
2020	Spasic <i>et al</i>	Clinical Text Data in Machine Learning: Systematic Review	(24)
2010	Stanfill <i>et al</i>	A Systematic Literature Review of Automated Clinical Coding and Classification Systems	(94)
2018	Velupillai <i>et al</i>	Using Clinical Natural Language Processing for Health Outcomes Research: Overview and Actionable Suggestions for Future Advances	(144)
2018	Wang <i>et al</i>	Clinical Information Extraction Applications: A Literature Review	(21)
2020	Wu <i>et al</i>	Deep Learning in Clinical Natural Language Processing: A Methodical Review	(23)

Questions?

bethany.percha@mssm.edu